



HAL
open science

NKNK: a New Essential Motif in the C-Terminal Domain of HIV-1 Group M Integrases

Marine Kanja, Pierre Cappy, Nicolas Levy, Oyindamola Oladosu, Sylvie Schmidt, Paola Rossolillo, Flore Winter, Romain Gasser, Christiane Moog, Marc Ruff, et al.

► **To cite this version:**

Marine Kanja, Pierre Cappy, Nicolas Levy, Oyindamola Oladosu, Sylvie Schmidt, et al.. NKNK: a New Essential Motif in the C-Terminal Domain of HIV-1 Group M Integrases. *Journal of Virology*, 2020, 94 (20), 10.1128/JVI.01035-20 . hal-02966338

HAL Id: hal-02966338

<https://hal.science/hal-02966338>

Submitted on 13 Oct 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **NKNK: a New Essential Motif in the C-Terminal Domain of HIV-1 Group M**

2 **Integrases**

3

4 Marine Kanja^{ab}, Pierre Cappy^{ab}, Nicolas Levy^c, Oyndamola Oladosu^c, Sylvie Schmidt^d, Paola
5 Rossolillo^{ab}, Flore Winter^{ab}, Romain Gasser^{ab}, Christiane Moog^d, Marc Ruff^c, Matteo Negroni^{ab#}, and
6 Daniela Lener^{ab#}

7

8 ^a *Université de Strasbourg, CNRS, Architecture et Réactivité de l'ARN, UPR9002, Strasbourg, France*

9 ^b *Interdisciplinary Thematic Institute (ITI) InnoVec, Strasbourg, France*

10 ^c *Chromatin Stability and DNA Mobility, Department of Structural Biology and Genomic, IGBMC,
11 Strasbourg University, CNRS, INSERM, Illkirch, France*

12 ^d *Molecular Immuno-Rheumatology Laboratory, UMR1109, FMTS, Université de Strasbourg, INSERM,
13 Institut de Virologie, Strasbourg, France*

14

15 Running head: HIV-1 M integrase new functional motif

16

17 # Corresponding authors: Daniela Lener, d.lener@ibmc-cnrs.unistra.fr, and Matteo Negroni,
18 m.negroni@ibmc-cnrs.unistra.fr. Mailing address: Institut de Biologie Moléculaire et Cellulaire, 2 allée
19 Konrad Roentgen, 67084 Strasbourg Cedex, France.

20

21 Abstract word count number: 219 (importance 116)

22 Main text word count number: 10 056

23 **ABSTRACT**

24 Using coevolution-network interference based on the comparison of two phylogenetically distantly
25 related isolates, one from the main group M and the other from the minor group O of HIV-1, we identify,
26 in the C-terminal domain (CTD) of integrase, a new functional motif constituted by four non-contiguous
27 amino acids (N₂₂₂K₂₄₀N₂₅₄K₂₇₃). Mutating the lysines abolishes integration through decreased 3'-
28 processing and inefficient nuclear import of reverse transcribed genomes. Solution of the crystal
29 structures of wt and mutated CTDs shows that the motif generates a positive surface potential that is
30 important for integration. The number of charges in the motif appears more crucial than their position
31 within the motif. Indeed, the positions of the K could be permuted or additional K could be inserted in
32 the motif, generally without affecting integration *per se*. Despite this potential genetic flexibility, the
33 NKNK arrangement is strictly conserved in natural sequences, indicative of an effective purifying
34 selection exerted at steps other than integration. Accordingly, reverse transcription was reduced even
35 in the mutants that retained wt integration levels, indicating that specifically the wt sequence is optimal
36 for carrying out the multiple functions integrase exerts. We propose that the existence of several amino
37 acids arrangements within the motif, with comparable efficiencies of integration *per se*, might have
38 constituted an asset for the acquisition of additional functions during viral evolution.

39 **IMPORTANCE** Intensive studies on HIV-1 have revealed its extraordinary ability to adapt to
40 environmental and immunological challenges, an ability that is also at the basis of antiviral treatments
41 escape. Here, by deconvoluting the different roles of the viral integrase in the various steps of the
42 infectious cycle, we report how the existence of alternative equally efficient structural arrangements for
43 carrying out one function opens on the possibility of adapting to the optimisation of further functionalities
44 exerted by the same protein. Such property provides an asset to increase the efficiency of the infectious
45 process. On the other hand, though, the identification of this new motif provides a potential target for
46 interfering simultaneously with multiple functions of the protein.

47

48 **Introduction**

49 Integration of reverse transcribed viral genomes into the genome of the infected cell is a peculiar feature
50 of the replication strategy of retroviruses, carried out by the viral enzyme integrase (IN) in a two-step
51 reaction. In HIV-1, after the achievement of DNA synthesis in the cytoplasm of the infected cell, it first
52 catalyses the removal of a conserved GT dinucleotide from the 3' ends of the viral DNA (3' processing),
53 leaving CA-OH 3' ends bound to the active site. Subsequently, once the viral DNA has been imported in
54 the nucleus, the reactive CA-OH-3' ends attack the cellular DNA leading to the generation of the provirus
55 (1, 2).

56 Besides this enzymatic function, HIV-1 IN is involved, through non catalytic activities, in several other
57 steps of the viral replication cycle. As a component of the Gag-Pol polyprotein precursor, it participates
58 in Gag-Pol dimerization, essential for the auto-activation of the viral protease and, consequently, for
59 viral particle maturation (3-5). During capsid morphogenesis, it is involved in the recruitment of the
60 genomic RNA inside the core of the viral particle (6). As a mature protein, it interacts with the viral
61 polymerase (reverse transcriptase, RT) to optimize reverse transcription of the viral genome (7-9).
62 Through the interaction with the cellular protein LEDGF/p75, it targets actively transcribed genes as
63 sites for integration (10). Finally, as a component of the pre-integration complex (PIC), IN is also involved
64 in nuclear import of the reverse transcription product, a peculiar feature of lentiviruses that allows the
65 infection of non-dividing cells.

66 This ability relies on the virus capacity to enter the nucleus *via* an active passage through the nuclear
67 pore complex (NPC) (11, 12). Several lines of evidence have indicated that the capsid (CA) protein is
68 crucial for nuclear entry (13, 14), through its interaction with several nucleoporins (Nups) forming the
69 NPC (Nup 358, Nup 153, Nup 98) (15-17) and with the transportin-3 (18, 19). Nevertheless, several
70 studies have indicated that IN has karyophilic properties. Namely, it contains a basic bipartite nuclear
71 localization signal (NLS) (20) as well as an atypical NLS (21), and also binds several cellular nuclear
72 import factors. Interactions of IN with importin α/β (22), importin 7 (23), importin $\alpha 3$ (24), Nup153 (25),
73 Nup62 (26) and transportin-3 (27-29) have been documented. Indeed, the mutation of amino acids,
74 mostly located in the C-terminal domain of the IN, responsible for binding to nuclear import factors,
75 results in non-infectious viruses impaired in nuclear import (23, 24, 28, 30).

76 The functional form of the HIV-1 integrase is made up of a dimer of dimers, which assemble in highly
77 ordered multimers of these tetramers (31, 32). Three domains, connected by flexible linkers, constitute
78 HIV-1 IN: the N-terminal domain (NTD), the catalytic core domain (CCD) and the C-terminal domain
79 (CTD) (33, 34). While the NTD is mostly involved in protein multimerization (35, 36), the CCD is mostly
80 responsible for catalysis, and for binding to the viral and cellular DNA as well as to the cellular cofactor
81 LEDGF (37-40). Finally, the CTD is involved in DNA binding during integration (41, 42), in protein
82 multimerization (35), in the interaction with the reverse transcriptase (7, 9) and in the recruitment of the
83 viral genomic RNA (gRNA) in the viral core (6). Overall, the intrinsic flexibility of the protein, the multiple
84 steps required to achieve integration, and the multimeric nature of the integration complex make the
85 involvement of the different parts of the protein in the various functions of the integrase very complex
86 and still not fully elucidated.

87 In addition, the multiple tasks that the IN must accomplish during the infectious cycle and the complexity
88 of its supramolecular structures are expected to impose functional constraints that ultimately may limit
89 its genetic diversity. Retention of functionality despite sequence variation strongly relies on covariation,
90 inside or outside the mutated protein. When an initial mutation negatively alters the protein functionality,
91 compensatory mutations can restore it, at least partially. Therefore, the sequences of homologous
92 proteins in different HIV variants are the result of independent evolution pathways, with independent
93 covariation networks specifically generated for each pathway. Chimeric genes between variants of a
94 given protein can perturb such networks and result in the production of non-functional proteins. This
95 information can then be exploited to probe the existence of functional motifs in proteins. For considerably
96 divergent viruses, as those derived from independent zoonotic transmissions, this approach can be
97 particularly powerful. This is the case for HIV-1 groups M and O that derive from simian viruses infecting
98 chimpanzees and gorillas, respectively. Here, we exploit the natural genetic diversity existing between
99 these groups to generate chimeric integrases. A detailed characterization of the individual amino acids
100 that differ in the non-functional chimeras has then led to the identification and functional characterization
101 of a new motif, in the CTD of HIV-1 group M integrase, essential for viral integration.

102

103

104 **Results**105 ***Analysis of intergroup M/O chimeras in the CTD of IN***

106 The functionality of the integrases studied in this work was evaluated following the protocol outlined in
107 Figure 1A-B and detailed in Materials and Methods. For this, we replaced the original RT and IN
108 sequences of the p8.91-MB (see Materials and Methods) by those of either one isolate of HIV-1 group
109 M, subtype A2, referred herein as "isolate A" or one isolate of HIV-1 group O, referred herein as "isolate
110 O". The resulting vectors were named RTA-INA (vRTA-INA) and (vRTO-INO), respectively. With these
111 vectors, we estimated the functionality of the integrases by measuring the efficiency of generation of
112 proviral DNAs. Since the number of proviral DNAs generated for each sample is dependent not only on
113 the levels of functionality of the integrase but also on the amount of total viral DNA generated after
114 reverse transcription, we estimated the amount of total viral DNA generated by each sample by qPCR
115 as described in Materials and Methods. In parallel, we measured the amount of proviral DNA generated
116 either by the puromycin assay or by the Alu qPCR assay as described in Materials and Methods
117 (Evaluation of integration by puromycin assay). The amount of proviral DNA divided by that of total viral
118 DNA provides an estimate of the efficiency of integration. Comparable efficiencies of integration were
119 measured with the two vectors, irrespective of whether the estimation was done using the puromycin
120 assay ($71 \pm 13\%$ and $72 \pm 24\%$ the level of the reference vector v8.91-MB respectively) or the Alu
121 qPCR assay ($70 \pm 14\%$ and $68 \pm 15\%$, respectively, Figure 1C). Throughout this work, the efficiency of
122 integration has always been evaluated by the puromycin-resistance assay normalised by the amount of
123 total DNA. Control vectors, in which the catalytic activity either of the integrase or of the reverse
124 transcriptase have been abolished, in vRTA-INA, by the introduction of the D116A mutation in IN or of
125 the D110N-D185N mutations in RT (43, 44), gave the expected results (Figure 1C).

126 We chose to probe the existence of functional motifs in the C-terminal domain of integrase because this
127 domain is involved in several non-catalytic functions of the protein. The C-terminal domain of INO used
128 in this study is 10 amino acids longer (212-298) than that of INA (212-288, Figure 2A). We constructed
129 three chimeras between isolates A and O, named after the position, in amino acids from the beginning
130 of the IN-coding region, where the sequence shifts from that of one isolate to that of the other (Figure
131 2B). Chimera A(1-212)-O(213-298) is constituted by INA with the entire CTD from INO; chimera A(1-

132 285)-O(286-298) is INA with the additional 10 amino acids of INO at the C-terminal end plus the two
133 most C-terminal different amino acids; finally, as the region between position 212 and 288 differs in 12
134 amino acids, chimera A(1-272)-O(273-298) was constructed in such a way as to split the 12 different
135 amino acids in two groups of 6.

136 We first performed western blots (Figure 2C) on viral particles to monitor the degree of proteolytic
137 processing of the Gag precursor (Pr55Gag), since incomplete processing would result in immature viral
138 particles, affecting infectivity. No significant differences in Pr55Gag were observed between isolate O
139 and chimerical constructs compared to isolate A (Figure 2D). We then evaluated the efficiency of reverse
140 transcription (measuring the amount of viral DNA produced by qPCR) and of integration (as described
141 above). Only chimera A(1-272)-O(273-298) exhibited significant defects in both reverse transcription
142 and integration (Figure 2E-F), suggesting that a covariation network, present between positions 212 and
143 285, was broken in this chimera. Since in these experiments the IN is expressed from p8.91-MB and
144 not from the genomic RNA, it can be ruled out that the phenotypes observed are due to an effect of the
145 mutations on the genomic RNA, as for example on the process of splicing, as it has been previously
146 described for some mutants of the C-terminal domain of IN (45).

147

148 **Characterization of IN CTD**

149 In order to evaluate the individual contribution of the 10 amino acids differing between positions 212 and
150 285 (Figure 2A), each residue in IN A was individually replaced by those of IN O and the ten-point
151 mutants were tested for processing of Pr55Gag, reverse transcription and generation of integrated
152 proviruses.

153 Except for mutant N254K, no significant difference in level of Pr55Gag proteolytic processing was
154 observed between mutants and parental vector A (Figure 3A). The effect on reverse transcription was
155 an overall reduction of efficiency for most of the mutants, with a residual efficiency between 45 and 90%
156 that of the parental vector A (Figure 3B). Concerning integration efficiency, instead, the majority of the
157 mutants did not show a significant decrease, except for mutants K240Q and K273Q for which integration
158 was dramatically impaired (Figure 3C). This suggests a specific implication of these two residues in the

159 integration process. When the two mutations were combined (K240Q/K273Q mutant), while the level of
160 reverse transcription remained above 40 % that of wt IN A, integration dropped to undetectable levels
161 (Figure 3D).

162 To discriminate between the role of the charge of K₂₄₀ and K₂₇₃ from that of their possible acetylation,
163 we replaced both residues by two R (K240R/K273R mutant). The level of integration of this mutant was
164 comparable to that of wt IN A (Figure 3E), indicating that the presence of a positive charge and not
165 acetylation at these positions was important for integration efficiency. However, these mutations
166 reduced by half both Pr55Gag proteolytic processing and reverse transcription (data available upon
167 request).

168 In both mutants showing a marked defect in integration (K240Q and K273Q), a K (positively charged
169 polar side chain) was replaced by a Q (non-charged polar side chain), the amino acid present in isolate
170 O at the corresponding positions. Conversely, in isolate O, two K are present in positions where a polar
171 non charged amino acid (N in both cases) is present in isolate A (positions 222 and 254, Figure 2A).
172 Therefore, in order to evaluate if also the two non-charged polar amino acids (N) present in isolate A at
173 positions 222 and 254 are essential, we replaced them by a non-polar amino acid like leucine (mutant
174 LKLK) and, in parallel, by a non-charged polar residue, Q (mutant QKQK). While in the LKLK mutant
175 the efficiency of integration dropped to almost undetectable levels, in the QKQK one it was comparable
176 to that of the wt enzyme, suggesting that the presence of a polar residue at these positions is essential
177 (Figure 4A). To understand whether the polar nature of the amino acid at positions 222 and 254 is
178 enough to retain functionality, the N were replaced by two threonine, which are polar but do not have
179 the amide group of asparagine. In this case (TKTK mutant) integration dropped to undetectable levels
180 (Figure 4A) indicating that not only the polarity is important but also the functional group carried by the
181 amino acid. Therefore, the biochemical features of all four residues identified are important.

182 Finally, we wondered whether the residues present at positions 222, 240, 254 and 273 could be
183 interchanged between isolates O and A. Therefore, we generated the quadruple mutant of isolate A
184 N222K/K240Q/N254K/K273Q (called KQKQ for simplicity). Remarkably, the integration efficiency of this
185 mutant was not significantly different from that of wt IN A (Figure 4B), indicating the existence of a
186 functional link between these four positions.

187 The alignment of HIV-1 IN sequences reveals a strong conservation of the amino acids N₂₂₂K₂₄₀N₂₅₄K₂₇₃
188 in group M (Figure 4C). To confirm the need for K₂₄₀ and K₂₇₃, observed in isolate A, also for other
189 isolates of group M, we introduced the K240Q/K273Q double mutation (NQNQ mutant) in integrases
190 from three other primary isolates of group M (Figure 4D). In all cases, a dramatic drop in integration was
191 observed with respect to the corresponding wt integrases, confirming the results obtained with isolate
192 A. The importance of the two K in the motif was therefore confirmed in isolates from the most widespread
193 HIV-1 group M subtypes in the epidemics, subtypes A, B, C, and CRF02 being responsible for 79% of
194 the HIV-1 infections worldwide (46).

195 The possibility of permuting the positions of the four amino acids at positions 222, 240, 254 and 273
196 indicates a functional relationship between these residues that can therefore be considered as a
197 functional motif that, based on the identity of the amino acids present at these positions in isolates of
198 group M, we refer to as the "NKNK" motif.

199

200 ***Importance of the lysines in the NKNK motif***

201 To understand to which extent the number and the positions of the K in the motif influence IN
202 functionality, we generated a series of mutants based on the replacement of the amino acids present in
203 isolate A by those of isolate O. We thus tested all possible variants (Figure 5A) containing either only
204 one K (four mutants, Figure 5B), two K (five mutants plus the wt, Figure 5C), three K (four mutants,
205 Figure 5D) or four K (one mutant, KKKK, Figure 5A) at any of the positions in the motif. The presence
206 of a single K led on average to a drop to 20% of integration with respect to wt IN A, whereas when two
207 or more lysines were present in the motif, levels of integration were close to those of wt IN A, ranging
208 from 75 to 137% (Figure 5A). The mutant with no K, where the motif sequence has been changed from
209 NKNK to NQNQ, confirmed the total loss of integration already observed with this mutant (see Figure
210 3D). Finally, from the analyses of the different mutants it appears that the presence of a K at the first
211 position of the motif (position 222) consistently leads to a higher level of integration in all classes of
212 mutants (those with 1, 2, or 3 K). Interestingly, though, position 222 has a N in the wt enzyme.

213 When considering individual mutants within the different classes, we observed a significant decrease in
214 functionality for all the mutants possessing only one K (Figure 5B). For the mutants containing two K,
215 three variants were at least as functional as wt IN A (NKNK in the figure), while two displayed a
216 significant reduction (Figure 5C). Finally, all mutants containing three K were at least as functional as
217 the parental IN A (Figure 5D). Remarkably, the results obtained with the mutants containing three or
218 four K indicate that the positively charged residues can replace the polar ones, while the reverse is not
219 the case, as shown by the mutants with none or only one K. Overall, these results indicate that at least
220 two K are required to have wt levels of integration, even if not all the positions in the motif are equivalent.
221 Instead, all mutants impacted reverse transcription with a reduction to 40-80% of the wt IN A (Figure
222 5E).

223

224 ***The NKNK motif in replication-competent viruses***

225 To confirm the observations made in the single infection cycle system, some mutants were then tested
226 in a replication-competent system using NL4.3 as primary virus. Mutants of the class containing two K
227 in the motif (the number of K found in circulating viruses) and with a marked phenotype were chosen for
228 this analysis. Besides the wt A sequence, we chose three mutants that either retained integration (KQKQ
229 and KQNK) or exhibited reduced integration (NQKK) (Figure 5C). To construct the four variants, we
230 replaced the sequence of NL4.3 CTD by that of isolate A, either wt or carrying the KQKQ, KQNK or
231 NQKK motifs (Figure 6A).

232 The infectivity of the virus carrying the whole CTD of INA instead of that of NL4.3 (called NL4.3 CTD A,
233 Figure 6A) was comparable to that of wt NL4.3 virus, set as reference, indicating that the replacement
234 of the whole CTD from NL4.3 by that of isolate A did not impact viral infectivity (Figure 6B). Regarding
235 the mutants, the results well recaptured the observations made with a single infection cycle (Figure 5C):
236 the infectivity was maintained for KQKQ and KQNK mutants while it was markedly decreased with NQKK
237 motif (Figure 6B).

238

239

240 **Role of the lysines of the motif in the integration process**

241 In order to characterise in which steps of the infectious cycle are involved the lysines of the NKNK motif,
242 we evaluated the effect of their mutation in two steps (other than reverse transcription) upstream the
243 integration of the pre-proviral DNA in the chromosomes of the host cell. In particular, by quantifying the
244 two LTR circles (2LTRc), we evaluated nuclear import and, by characterizing the LTR-LTR junctions of
245 2LTRc, the efficiency of 3' processing, which takes place in the cytoplasm, before nuclear import.

246 2LTRc are exclusively formed in the nucleus and are, therefore, useful markers for nuclear import of the
247 reverse transcribed genomes (47). They are generated when the full-length reverse transcription
248 products are not used as substrate for integration. If a mutant is defective in catalysis but carries out
249 nuclear import efficiently (as mutant D116A), 2LTRc should accumulate with respect to a wt IN. Instead,
250 if the mutant is also impaired in nuclear import, 2LTRc will either not increase with respect to the wt IN
251 or increase but more modestly than for D116A.

252 Hence, to monitor nuclear import, we measured the amount of (2LTRc) in wt IN A, in mutants containing
253 either no K (NQNQ) or only one (either K₂₇₃, NQNK mutant, or K₂₄₀, NKNQ mutant). IN A D116A mutant
254 was used as a control. This mutant being totally inactive for integration was considered to produce the
255 highest accumulation of 2LTRc, set at 100%. As expected, the level of 2LTRc found with wt IN A, which
256 efficiently imports and integrates the reverse transcribed genome, was significantly lower (25%) than
257 that of the D116A mutant. As shown in Figure 7A, despite their inability to generate proviral DNA, the
258 mutants had levels of 2LTRc significantly lower than IN A D116A, indicative of a defect in nuclear import.

259 To estimate the efficiency of nuclear import in the mutants, we estimated the level of 2LTRc (Table 1,
260 line 2, "theoretical level"), that could be obtained if no defect in nuclear import was present. We then
261 calculated the efficiency of nuclear import as the ratio between the level of 2LTRc observed
262 experimentally (Table 1, line 3) and the theoretical one. If no defects in nuclear import are present, ratios
263 should be around 1, while defects in nuclear import would yield ratios <1. The ratios found for the three
264 mutants were in the 0.31-0.35 range (Table 1, line 4), indicative of a reduction of nuclear import to
265 approximately 1/3 that of the wt enzyme. Therefore, the defects in nuclear import contribute to the
266 decrease in integration found with these mutants, but cannot alone account for the low levels observed,

267 particularly in NQNQ and NQNK mutants for which integration was almost undetectable (Table 1, line
268 1).

269 The efficiency of 3' processing carried out by IN was then analysed by quantifying the different LTR-
270 LTR junctions in the 2LTRc. The 2LTRc found in the nucleus are generated from DNAs carrying either
271 unprocessed or processed 3' ends. In the first case, the 2LTRc will present "perfect junctions" (PJ) while
272 in the second the junctions will be "imperfect". A high ratio of PJ/2LTRc is therefore indicative of
273 inefficient 3' processing. We found that mutating both K (mutant NQNQ) or only K₂₄₀ (mutant NQNK) led
274 to results not significantly different from those obtained with the IN A D116A catalytic mutant (Figure
275 7B), indicative of a marked defect in 3' processing. Mutating K₂₇₃ (mutant NKNQ), instead, did not affect
276 the process, with PJ/2LTRc values comparable to those of the wt enzyme.

277 To evaluate the contribution of defects in 3' processing to the decreased integration efficiency observed
278 with the various mutants, we first estimated the maximum diminution of PJ/2LTRc ratio observed for a
279 fully competent enzyme (wt IN A). The PJ/2LTRc ratio for wt IN A was 0.54 that of IN A D116A (Table
280 1, line 5), corresponding to a reduction of 46 % due to 3' processing (Table 1, line 6). The ratio PJ/2LTRc
281 with respect to IN A D116A was then calculated for each mutant and the resulting value was divided by
282 0.46, obtaining an estimate of the efficiency of 3' processing relative to that observed for wt IN A (Table
283 1, line 7). 3' processing of NQNQ and NQNK mutants was dramatically reduced, 15 % and 28 % that of
284 wt IN A, respectively. Mutating K₂₇₃, instead, only decreased 3' processing to 74 % that of wt IN A.

285 Finally, in order to understand if these two types of defects (nuclear import and 3' processing) were
286 sufficient to explain the integration defects observed with the various mutants, we combined the effect
287 of these defects (Table 1, line 8). The values obtained account remarkably well for the efficiencies of
288 integration observed (lines 1 and 8 of Table 1) indicating that the decrease observed when mutating the
289 K of the NKNK motif, once normalized for the differences observed in the amount of viral DNA produced,
290 is essentially due to alterations in these two processes.

291

292

293

294 **Structural analysis of wt and mutant integrase C-terminal domains**

295 To understand the structural bases for the functional differences observed in the NKNK motif mutants,
296 the crystal structures of the C-terminal domain (IN CTD, 220-270) of wild type IN A and of the reference
297 strain NL4.3 were solved at 2.2 Å and 1.3 Å of resolution, respectively. For both structures, K₂₇₃ was not
298 included as it is in a disordered region of IN. For all crystal forms we observed a strong packing
299 interaction through the His-tag coordinating a Nickel ion (data available upon request). The structures
300 had the same topology, consisting in a five-stranded β-barrel (data available upon request). The region
301 encompassing the positions of the motif (Figure 8A-B) generates a surface endowed with a positive
302 potential (circled in yellow in Figure 8C-D), suggesting that this feature could be important for the
303 functionality of the IN. In this case, it is expected that inserting additional lysines in the motif (as for the
304 mutants containing three or four lysines) would retain functionality and, conversely, removing the K
305 (mutants with one K or no K) would affect it. This is what we observed in Figure 5. Nevertheless, the
306 correlation between surface potential and functionality is less clear for the mutants where the number
307 of K in the motif (two) is not altered but their positions are permuted with polar amino acids.

308 To clarify this point, we solved, at a 2.0 Å resolution the crystal structure of the CTD of the NQKK mutant,
309 which was the one displaying the most dramatic drop in integration among the mutants possessing two
310 K (Figure 5C). The NQKK CTD crystallized in a different space group and had three chains in the
311 asymmetric unit. The superposition of the five structures corresponding to NL4.3 CTD, to A CTD and to
312 the three molecules in the asymmetric unit for the NQKK CTD (chains A, B and C) did not show
313 significant differences in the main chain fold (Root-mean-square deviation of atomic positions, RMSD:
314 NL4.3 CTD vs A CTD = 0.395 Å, A CTD vs NQKK CTD ABC = 0.653 Å, NQKK CTD chain A vs chain B
315 vs chain C = 0.558 Å). Interestingly, the positive surface electrostatic potential observed for the wt
316 enzyme was markedly perturbed in the NQKK mutant (yellow circle in Figure 8E-F), a change that could
317 well account for the decrease of functionality of the NQKK mutant integrase.

318 To further analyse the impact of the mutations on the structure, we defined the regions which are
319 naturally disordered (Intrinsically Disordered Regions, IDRs) by the superposition of the three molecules
320 in the asymmetric unit of the NQKK CTD structure. We assume that the change in the RMSD obtained
321 among the three molecules represents the natural IDRs. For the main chain, disordered regions with

322 high RMSD are 228-232 and 243-248 (data available upon request). These same regions are found to
323 be similarly disordered when comparing the main chain RMSD of the C-terminal domain of IN A to that
324 of IN A NQKK chains A, B and C (data available upon request), indicating that the mutations have no
325 effect on the C-alpha backbone fold of IN CTD.

326 The three structures were then analysed from the standpoint of the arrangement of the side chains.
327 Calculating side chains RMSD, disordered portions were found to correspond to regions 222-225, 228-
328 232 and 243-248 (data available upon request). The comparison of A CTD and NL4.3 CTD, which differ
329 between positions 220-270 only by a single amino acid change (V234 in NL4.3 replaced by I234 in A),
330 expectedly, did not reveal significant differences between the two structures. Instead, when comparing
331 A CTD to NQKK CTD A, B and C side chain deviation, we observed a difference in the structure of
332 region 235-237 (data available upon request). This difference appears to be due to the N254K mutation
333 that induces a displacement of the side chain of lysine 236 (white arrows in Figure 8G). This is likely
334 due to a repulsive interaction between the side chains of the two lysines, resulting in a perturbation of
335 the structure in the 235-237 region.

336 To evaluate the importance of charge configurations in the context of the C-terminal domain of retroviral
337 integrases, we performed an analysis of the electrostatic charge surface potential for other lentiviruses
338 as well as for other retroviruses. Integrases C-terminal domain structures are available for HIV-1 A2,
339 PDB 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V (48); MVV,
340 PDB 5LLJ (49); RSV, 1COM (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV, PDB 4E7I
341 (53). First, we superposed the available structures. The superposition shows that they share a common
342 fold (Figure 9A) as well as a low Root Mean Square Deviation on secondary structure backbone despite
343 levels of sequence identity for some cases as low as 10% (data available upon request). A structure-
344 based sequence alignment has then been performed (Figure 9B). Surprisingly, despite a low overall
345 sequence identity (10 – 20 %) for some integrases, several regions have a strong local sequence
346 similarity (red and yellow background in Figure 9B) while no conservation is observed at positions 222,
347 240 and 254 among lentiviruses nor retroviruses. However, when we compared the electrostatic surface
348 potentials of all structures (Figure 9, Panels C-L), we could define two general retroviral classes. A first
349 class represented by lentiviruses where the surface corresponding to the one delimited by the NKNK

350 motif in HIV-1 M is basic and a second class represented by the other retroviruses tested
351 (orthoretroviruses α , β , γ , and spumaretroviruses) where this surface is acidic or neutral.

352

353

354 **Discussion**

355 Here, by performing a systematic comparison between the non-conserved amino acids in the CTD of
356 the HIV-1 group M and group O integrases, we identify, in group M, a highly conserved motif that is
357 essential for integration. The motif is constituted of two asparagines and two lysines (N₂₂₂K₂₄₀N₂₅₄K₂₇₃)
358 all required for the generation of proviral DNA. In particular, when the K were mutated, integration was
359 abolished due to the cumulative effects of decreased 3' processing and nuclear import of the reverse
360 transcription products (Table 1). Replacing the K by R did not affect integration (Figure 3E), suggesting
361 that the essential feature of the K is their positive charge. Importantly, the positions of the two K of the
362 motif could be permuted without affecting the functionality of the integrase in most cases (Figure 5).

363 A potential explanation for the retention of functionality when permutating the positions of the K across
364 the motif comes from the structural data on the CTD obtained in this work. We have solved the crystal
365 structure of the CTD of the wt IN A used in this study as well as that of the K240Q-N254K mutant
366 (referred as NQKK in the result section). In the structure of the wt enzyme, the residues constituting the
367 motif (except K₂₇₃ which is part of an unresolved region) generate a positively charged surface (Figure
368 8A-D). This positive electrostatic potential surface is absent, instead, in the NQKK mutant (Figure 8E-
369 F) which, despite the presence of two K, displays a drastic reduction of integration efficiency. These
370 results, combined to the tests of functionality of the different mutants, suggest that the relevant
371 parameter is the presence of a positive charge across this surface. Charged residues have a strong
372 effect on the surface potential. The nature of the amino acid side chain (charged, polar, non-polar) on
373 the surface of the protein defines the surface potential. Charged and polar groups, through forming ion
374 pairs, hydrogen bonds, and other electrostatic interactions, impart important properties to proteins.
375 Modulation of the charges on the amino acids, e.g. by pH and by post-translational modifications, have
376 significant effects on protein – protein and protein – nucleic acid interactions (54). In addition to residues

377 carrying net charges, also polar residues have significant partial charges and can form hydrogen bonds
378 and other specific electrostatic interactions among themselves and with charged residues (55). In the
379 case of the present study, the possible contribution of these mechanisms to the functionality of the
380 integrase could be reflected by the loss of functionality observed by replacing the N, which carries an
381 amide side chain ($-\text{CONH}_2$), by either a non-polar amino acid (L) or by a polar one (T) but carrying a
382 hydroxyl side chain ($-\text{OH}$). The analysis of the electrostatic surface potential for the integrases C-
383 terminal domain of the retroviruses for which this is known showed that the topology of the structure is
384 maintained (Figure 9A) and the analysis of the surface electrostatic potential splits the viruses studied
385 in two classes. One, constituted by lentiviruses, for which the surface delimited by the NKNK motif of
386 HIV-1 M contains basic charges (in some cases brought by amino acids non corresponding to those of
387 the NKNK motif of HIV-1 M). The second class including the other orthoretroviruses studied (α , β , γ)
388 and spumaretroviruses where this surface contains acidic and neutral regions. This presence of basic
389 regions, specifically in lentiviruses, could contribute to some specificity of lentiviral biology as to increase
390 the efficiency of infection of quiescent cells.

391 The importance for protein functionality of charge configurations and clusters in their three-dimensional
392 structures has been underlined by several studies (56-59). Charge permutations have been used in the
393 NC region of the Gag protein for the Mason-Pfizer Monkey Virus (60). This basic region could be
394 replaced with nonspecific sequences containing basic amino acid residues, without altering its
395 functionality while mutants with neutral or negatively charged residues showed a large drop in viral
396 infectivity in single round experiments. Moreover, a mutant exhibiting an increased net charge of the
397 basic region, was 30% more infectious than the wild type. Also, in our study, increasing the positive
398 charge of the NKNK motif of HIV-1 IN by introducing a third K leads to a slight increase in integration
399 with respect to wt IN (Figure 5A and D).

400 As retention of IN functionality relies on the electrostatic surface potential rather than on the specific
401 positions of the positively charged amino acid, we infer that this region is probably involved in the
402 interaction with a partner carrying a repetitive negatively charged biochemical motif, as the phosphates
403 of the nucleic acids backbone. Alternatively, the partner could be a disordered region of a protein that
404 can rearrange to preserve the interaction when the positions of the positive charges are permuted
405 across the surface of the NKNK motif. Indeed, the molecular recognition between charged surfaces and

406 flexible macromolecules like DNA, RNA and intrinsically disordered protein regions has been observed
407 for the Prototype Foamy Virus and Rous Sarcoma Virus Gag precursors (61, 62), for EBNA proteins of
408 the Epstein-Barr virus (63), UL34 protein of the Herpes Simplex Virus (64) as well as for cellular proteins
409 like APOBEC3G (65). Moreover, the presence of asparagines in the motif, which we show are required
410 for integrase functionality, could contribute to the interaction with the nucleic acid or with a protein
411 partner through hydrogen bonds with the bases (54) or with polar amino acids (55, 66, 67) , respectively.
412 The analysis of the motif in the context of the well-characterised structure of the intasome (31),
413 mimicking a post-integration desoxyribonucleic complex, indicates that the residues forming the
414 electrostatic surface point toward the solvent, at the exterior of the structure (Figure 10). This is coherent
415 with the observation that mutating the motif does not affect late steps of the integration process, but
416 rather earlier ones as 3' processing and nuclear import.

417 We show that the NKNK motif of the CTD is involved in 3' processing and nuclear import of the reverse
418 transcription product. Indeed, the removal of the K impacts both processes and when combined, these
419 effects are sufficient to account for the drop of infectivity to the undetectable levels observed in the
420 absence of K in the motif (Table 1). Since mutating K₂₄₀ has a strong impact on both 3' processing and
421 nuclear import, while K₂₇₃ appears to be predominantly involved in nuclear import, it is possible that the
422 involvement of the motif in these two processes implicates structurally distinct functional complexes.
423 This is the first finding of an implication of the HIV-1 IN CTD in 3' processing. So far, only the involvement
424 of the catalytic domain had been demonstrated (39, 68, 69). Since it has been shown that different
425 oligomerization states of IN influence specifically the ability to carry out 3' processing or strand transfer
426 (70), it is possible that the electrostatic surface formed by the NKNK motif help stabilize the oligomeric
427 state that allows 3' processing.

428 Concerning nuclear import, it is known that HIV-1 IN binds several cellular nuclear import factors through
429 basic amino acids of the CTD, and that abolishing these interactions leads to non-infectious viruses
430 displaying a severe defect in nuclear import. Here, we extend the regions of IN involved in this process
431 by describing the need for a new motif, although it cannot be discriminated whether its involvement is
432 direct or mediated by the interaction with a partner protein with karyophilic properties.

433 Some of the residues constituting the N₂₂₂K₂₄₀N₂₅₄K₂₇₃ motif have been previously characterized
434 showing their implication in different steps of the infectious cycle. One is the involvement in reverse
435 transcription. The integrase CTD interacts with the reverse transcriptase to improve DNA synthesis (71,
436 72). In one study, the double mutation K240A/K244E caused a decrease in reverse transcription to
437 around 20 % the levels of the wt enzyme (72) while the K244E mutation alone caused a reduction of 40
438 % of RT efficiency (73), suggesting that K₂₄₀ also contributes to reverse transcription. The decrease we
439 observed when mutating K₂₄₀ alone, to around 45 % of wt activity, is consistent with this view. The
440 characterisation by NMR of the RT-binding surface in the IN CTD, obtained using the CTD₂₂₀₋₂₇₀, shows
441 that it is made up of 9 residues (amino acids 231-258 among which K₂₄₄) that strongly interact both with
442 the RT alone (9) and with the RT/DNA complex (74). When the interaction involves the complex, this
443 surface includes 5 additional amino acids (74). Among these additional residues are N₂₂₂ and K₂₄₀, which
444 are located at one edge of the surface. It is therefore possible that the nature of the residues at positions
445 222 and 240 affects the interaction between the CTD and RT/DNA complex.

446 Concerning K₂₇₃, contradictory results have been obtained for reverse transcription of viruses harbouring
447 integrases with sequential C-terminal deletions (IN₁₋₂₇₀ and IN₁₋₂₇₂) (75, 76). Furthermore, for reverse
448 transcription to occur, the genomic RNA must be encapsidated in the core of the viral particle. In this
449 sense, it has been recently shown that mutating K₂₇₃ together with R₂₆₉ (R269A/K273A mutant) impairs
450 encapsidation of the genomic RNA (6). As expected, reverse transcription in the double mutant was
451 almost abolished. Here, mutating K₂₇₃ to Q led only to a reduction of reverse transcription of 30% (Figure
452 3B), suggesting that mutating K₂₇₃ alone is not sufficient to affect genomic incorporation into the viral
453 capsid, at least in the majority of the particles. Supporting this, an earlier study (77) showed that the
454 K273A single mutation did not affect viral replication in Jurkat cells, indicating that is the specific
455 combination of R269A/K273A mutations to be responsible for the impairment of the genome
456 encapsidation.

457 Finally, acetylation of K₂₇₃, has been previously proposed to be important for different steps of the
458 infectious cycle (78, 79). In those studies, though, the role of acetylation of K₂₇₃ was assessed by
459 simultaneously replacing K₂₆₄, K₂₆₆, and K₂₇₃, thereby not allowing to conclude on the specific
460 contribution of K₂₇₃. Here, hampering acetylation but preserving the positive charge by the K273R
461 substitution did not affect integration, indicating that the possible acetylation of K₂₇₃ had no effect on

462 integration. This observation is in line with the observation by Topper and co-workers that
463 posttranslational acetylation of the integrase CTD is dispensable for viral replication (80). Altogether,
464 the data available in the literature regarding K₂₄₀ and K₂₇₃ indicate that the effects we observed in this
465 study cannot be due to any of the already known properties of the residues of the motif.

466 The NKNK motif is strictly conserved in natural sequences of HIV-1 group M. However, we show that
467 various variants of the NKNK motif display levels of integration efficiency equivalent to the wt enzyme
468 and could therefore in principle be found in the epidemics. Their absence is indicative of purifying
469 selection occurring *in vivo*, likely exerted at a step different from integration. One possibility is the
470 implication of IN in reverse transcription, which is, in all variants, less efficient than with the NKNK
471 sequence. The existence of several alternative sequence arrangements, in the motif, with comparable
472 efficiencies of integration might therefore have constituted an asset for optimizing the acquisition of
473 additional functions, such as promoting reverse transcription.

474

475 **Materials and Methods**

476 **Plasmids and molecular cloning**

477 p8.91-MB was constructed by engineering one *MluI* and one *BspEI* restriction sites respectively 18 nt
478 downstream the 5' and 21 nt upstream the 3' of the RT-coding sequence of the pCMVΔR8.91 (81). The
479 insertion of the two restriction sites led to three amino acids changes in the RT (E6T, T7R and A554S).
480 These modifications only slightly affected the efficiency of generation of puromycin-resistant clones (see
481 below) upon transduction with the resulting viral vector (v8.91-MB) since, in three independent
482 experiments, the number of clones obtained with p8.91-MB was 80% ± 6% of that obtained using p8.91.
483 The p8.91-MB was employed throughout the study as positive control and was used to insert the various
484 variants of RT and IN tested. Together with the *SalI* site, present in the p8.91 48 bp downstream the
485 stop codon of *pol* gene CDS, the *MluI* and *BspEI* sites define two exchangeable cassettes: one
486 encompassing the RT coding sequence (*MluI*-RT-*BspEI*, 1680 bp) and one encompassing the IN-coding
487 sequence (*BspEI*-IN-*SalI*, 940 bp). These cassettes were used to insert the various sequences of RT
488 and IN used in the study. The plasmid used to produce the genomic RNA of the viral vectors was a

489 modified version of pSDY, previously described (82), hereafter called pSRP (for pSDY-nRFP-Puro).
490 This variant was obtained by introducing two modifications to the original pSDY-dCK-Puro plasmid (82).
491 The first one was the replacement of the sequence encoding the human deoxycytidine kinase by a
492 cassette containing the RFP fused with the N-terminal 124 amino acids of human histone H2B, which
493 directs the RFP to the nucleus. The RFP was used to monitor the efficiency of transfection by
494 fluorescence microscopy. The second modification was the replacement of the HIV-1 U3 sequence in
495 the 5' LTR by that of the U3 of the Rous sarcoma virus. For the generation of qPCR standard curves,
496 two plasmids were constructed: one, called pJet-1LTR, for the detection of early and late reverse
497 transcription products, was obtained by inserting the sequence encompassing the LTR and the Psi
498 region from pSDY (82) in the pJET plasmid with the CloneJET PCR Cloning Kit (Thermo Scientific, MA,
499 USA); the second, pGenuine2LTR, has been obtained by inserting a fragment of 290 bp corresponding
500 to the unprocessed junction of U5/U3 (CAGT/ACTG being the sequence of the junction 5' to 3') into the
501 pEX-A2 plasmid (Eurofins Genomics, Luxembourg). For the study with replication-competent viruses
502 we used the pNL4.3 plasmid (83) that was obtained from the NIH AIDS Research and Reference Reagent
503 Program, #114 (GeneBank accession #AF324493). We replaced in this plasmid the coding sequence
504 of NL4.3 IN CTD with those of wt and mutants INA CTD, as described in Results. Chimerical integrases
505 between primary isolates from HIV-1 group M subtype A2 and HIV-1 group O RBF206, as well as mutant
506 integrases, were constructed through overlap extension PCR as previously described for the envelope
507 gene (84).

508

509 **Cells**

510 HEK-293T cells were obtained from the American Type Culture Collection (ATCC). P4-CCR5 reporter
511 cells are HeLa CD4⁺ CXCR4⁺ CCR5⁺ carrying the LacZ gene under the control of the HIV-1 LTR
512 promoter (85). TZM-bl cells are a HeLa cell clone genetically engineered to express CD4, CXCR4, and
513 CCR5 and containing the Tat-responsive reporter gene for the firefly luciferase under the control of the
514 HIV-1 long terminal repeat (86). HEK-293T, P4-CCR5 and TZM-bl cells were grown in Dulbecco's
515 Modified Eagle's Medium (DMEM, Thermo Fisher, MA, USA) supplemented with 10% foetal calf serum
516 and 100 U/ml penicillin-100 mg/ml streptomycin (Thermo Fisher, MA, USA) at 37°C in 5 % CO₂. CEM-

517 SS cells are human T4-lymphoblastoid cells (87-89) and were grown in Roswell Park Memorial Institute
518 medium (RPMI) supplemented with 10% foetal calf serum and 100 U/ml penicillin-100 mg/ml
519 streptomycin (Thermo Fisher, MA, USA) at 37°C in 5 % CO₂.

520

521 **Viral strains**

522 The following primary isolates were used for this study: from HIV-1 group M, one from subtype A2
523 (GenBank accession #AF286237, named hereafter "isolate A"), one from subtype C (GenBank
524 accession #AF286224, hereafter named "isolate C"), one from CRF02_AG (GenBank accession
525 #MH351678), one from subtype B (isolate AiHo GenBank accession #MH351679, hereafter named
526 isolate B); from HIV-1 group O the primary isolate RBF 206 (GenBank accession #KU168298, hereafter
527 named "isolate O"). Isolates #AF286237, #AF286224 and #MH351678 were obtained from the NIH
528 AIDS Research and Reference Reagent Program; isolates #MH351678, #MH351679 and #KU168298
529 were kindly provided by J.C. Plantier (CHU Rouen, France).

530

531 **Sequence alignments**

532 We used 3366 HIV-1 sequences for alignment. HIV-1 group M sequences were downloaded from the
533 Los Alamos National Laboratory (LANL) HIV sequence database and correspond to the different HIV-1
534 group M pure subtypes: A (249 sequences), B (2450 sequences), C (450 sequences), D (121
535 sequences), G (80 sequences), H (8 sequences), J (6 sequences), K (2 sequences). We also aligned
536 49 HIV-1 group O sequences, using 26 sequences from the LANL database and the 23 sequences
537 obtained through collaboration with the Virology Unit associated to the French National HIV Reference
538 Center (Pr. J.C. Plantier). Sequence alignments were performed with CLC sequence viewer 8. The
539 sequence logo of positions 222, 240, 254 and 273 in HIV-1 group M IN was obtained with an alignment
540 of 3366 sequences of the IN CTD using WebLogo (<http://weblogo.berkeley.edu/logo.cgi>).

541

542

543 **Generation of pseudotyped viral vectors**

544 Pseudotyped lentiviral vectors were produced by co-transfection of HEK 293T cells with pHCMV-G (90)
545 encoding the VSV-G envelope protein, pSRP and p8.91-MB based plasmids with the polyethylenimine
546 method following the manufacturer's instructions (PEI, MW 25000, linear; Polysciences, Warrington,
547 PA, USA). HEK 293T were seeded at 5×10^6 per 100-mm diameter dish and transfected 16-20h later.
548 The medium was replaced 6h after transfection, and the vectors were recovered from the supernatant
549 72h later, filtered on 0.45 μm filters and the amount of p24 (CA) was quantified by ELISA (Fujirebo
550 Europe, Belgium).

551

552 **Western blot**

553 Western blot analysis was carried out on virions to assess the proteolytic processing of the Pr55Gag
554 polyprotein. 1.5 mL of viral supernatant was centrifuged through 20 % sucrose, and the virion pellet was
555 lysed in Laemmli buffer 1.5X. Viral proteins were separated on a Criterion™ TGX Strain-Free 4-15 %
556 gradient gel (Biorad, CA, USA) (TGS: Tris Base 0,025 M/Glycine 0,192 M/SDS 0,1 %, 150V, 45 min),
557 blotted on a PVDF membrane (TGS/Ethanol 10 %, 200 mA, 1.5h) and probed with a mouse monoclonal
558 anti-CA antibody (NIH AIDS Reagent Program, #3537) to detect the viral capsid, the Pr55Gag
559 unprocessed polyprotein and CA-containing proteolytic intermediates. An anti-mouse HRP-conjugated
560 secondary antibody was used to probe the membrane previously incubated with anti-CA. Membranes
561 were incubated with ECL reagent (Thermo Fisher, MA, USA) and WB were imaged on a Biorad
562 Chemidoc Touch and analysed with the Biorad Image Lab software.

563

564 **Evaluation of reverse transcription by qPCR**

565 The viral vectors were treated with 200 U/ml of Benzonase nuclease (Sigma-Aldrich, MO, USA) in the
566 presence of 1 mM MgCl_2 for 1h at 37°C to remove non-internalized DNA. The vectors (200 ng of p24)
567 were then used to transduce 0.5×10^6 HEK 293T cells by spinoculation for 2h at 32°C, 800 rcf, with
568 8 $\mu\text{g}/\text{mL}$ polybrene (Sigma-Aldrich, MO, USA). After 2h, the supernatant was removed, cells were

569 resuspended in 2 mL of DMEM and plated in 6-well plates. After 30h, cells were trypsinised and pelleted.
 570 Total DNA was extracted with UltraClean® Tissue & Cells DNA Isolation Kit (Ozyme, France). A duplex
 571 qPCR assay (see Table 3 for primers) was used to quantify early and late reverse transcription products
 572 by detecting the R-U5 and U5-Psi junctions, respectively, and another qPCR (Table 3) to normalise for
 573 the quantity of cells employed in the assay (detection of β -actin exon 6 genomic DNA; International DNA
 574 Technologies -IDT- Belgium). All primers and probes were synthesised by IDT. The qPCR assays were
 575 designed with the Taqman® hydrolysis probe technology using the IDT Primers and Probes design
 576 software (IDT), with dual quencher probes (one internal ZEN™ quencher and one 3' Iowa Black™ FQ
 577 quencher) (Table 3). qPCRs were performed with the iTaq Universal Probes Supermix (Biorad, CA,
 578 USA) on a CFX96 (Biorad, CA, USA) thermal cycler with the following cycling conditions: initial Taq
 579 activation 3 min, 95°C followed by [denaturation 10 sec/95°C; elongation 20 sec/55°C] x 40 cycles.
 580 Standard curves and analysis were carried out with the CFX Manager (Biorad, CA, USA). DNA copy
 581 number was determined using a standard curve prepared with serial dilutions of the reference plasmids
 582 pJet-1LTR and of a known number of HEK 293T cells.

583

584 **Evaluation of integration by puromycin assay**

585 Half a million of HEK 293T cells were transduced with a volume of viral vectors corresponding to 0.2 ng
 586 of p24, by spinoculation 2h at 32°C, 800 rcf, with 8 μ g/mL polybrene (Sigma-Aldrich, MO, USA). After 2h
 587 the supernatant was removed, cells were resuspended in 7 mL of DMEM and plated in 100 mm diameter
 588 plates. After 30h, puromycin was added at a final concentration of 0.6 μ g/mL, clones were allowed to
 589 grow for 10 to 12 days and then counted. However, the number of clones depends on two parameters:
 590 the efficiency of integration and the amount of pre-proviral DNAs available for integration (which
 591 depends on the efficiency of reverse transcription). Therefore, we normalized the number of clones
 592 observed by the amount of viral DNA generated by reverse transcription (estimated by qPCR) to
 593 extrapolate the efficiency of integration. The percentage of integration efficiency for sample X with
 594 respect to the control C is, thus, given by $(p_x/r_x)/(p_c/r_c) \times 100$, where r_x and r_c are the amounts of late
 595 reverse transcription products, estimated by qPCR, in sample X and in control C, respectively, and p_x
 596 and p_c the number of puromycin-resistant clones in sample X and in control C, respectively.

597

598 Evaluation of integration by Alu qPCR

599 Equal amounts of total DNA extracted from transduced cells (as deduced by qPCR of β -actin exon 6
600 genomic DNA, see above) were used for the Alu PCR assay, as previously described (91). Two
601 subsequent amplification were performed. The first one, 95°C for 3 min, [95°C for 30 sec, 55°C for 30
602 sec, 72°C for 3 min30s] x15, 72°C for 7 min, using the Alu-forward primer and the Psi reverse primer,
603 allowed to amplify Alu-LTR fragments (Table 3). Samples were then diluted to 1:10 and 2 μ L were used
604 for the second amplification to detect the viral LTR, as described above for the detection of the R-U5
605 junction. The percentage of integration efficiency for sample X with respect to the control C is given by
606 $(a_x/r_x)/(a_c/r_c) \times 100$, where r_x and r_c are the amounts of late reverse transcription products, estimated by
607 qPCR, in sample X and in control C, respectively, and a_x and a_c the amounts of DNA estimated by the
608 second amplification of the Alu qPCR assay in sample X and in control C, respectively.

609

610 Quantification of two LTR circles and of circles with perfect junctions

611 Non-internalised DNA was removed by treatment with Benzonase nuclease as for the qPCR assay and
612 0.5×10^6 HEK 293T cells were transduced with a volume of viral vectors corresponding to 1 μ g of p24
613 by spinoculation, as described above. After 30h, cells were trypsinised and pelleted. Total DNA was
614 extracted with UltraClean® Tissue & Cells DNA Isolation Kit. Late reverse transcription products
615 (detection of the U5-Psi junction) were quantified as described above and two qPCR assays were used
616 to quantify the 2LTRc and the quantity of 2LTR circles with a perfect palindromic junction, with a primer
617 overlapping the 2LTRc junction, as previously described (92). The qPCR assays were designed and the
618 primers and probes (Table 4) synthesised as described above. qPCRs were performed as described
619 above. Standard curves and analysis were carried out with the CFX Manager (Biorad, CA, USA). Copy
620 numbers of the different forms of viral DNA were determined with respect to a standard curve prepared
621 by serial dilutions of the pGenuine2LTR plasmid. The amount of 2LTRc for sample X is normalised by
622 the total amount of viral DNA (detection of the late reverse transcription products by quantifying the U5-
623 Psi junction), then it is expressed as a percentage of the amount detected for the control INA D116A

624 (indicated with D), thus giving $(2LTR_{Cx}/r_x)/(2LTR_{CD}/r_D) \times 100$, where r_x and r_D are the amounts of late
 625 reverse transcription products, in sample X and in control D, respectively, and $2LTR_{Cx}$ and $2LTR_{CD}$ the
 626 amount of 2LTR circles in sample X and in control D, respectively.

627

628 **Calculation of the efficiency of nuclear import and of 3' processing**

629 The efficiency of nuclear import was estimated as follows. The level of 2LTRc found with wt IN A was
 630 0.2 with respect to that found with D116A (data from Figure 7A). The diminution observed with wt IN A
 631 with respect to D116A (which was considered to produce the maximum amount of 2LTRc and was
 632 therefore set at 1) was therefore 0.8 (given by $1 - 0.2$). Since the diminution of 2LTRc is proportional to
 633 the efficiency of integration, for example a mutant integrating with an efficiency 0.3 that of wt IN A is
 634 expected to reduce the amount of 2LTRc by $0.8 \times 0.3 = 0.24$. The amount of 2LTRc expected for that
 635 mutant would therefore be given by $1 - 0.24 = 0.76$. Similarly, a mutant integrating with a higher efficiency
 636 (for example 0.9 that of wt IN A) is expected to give $1 - (0.8 \times 0.9) = 0.28$ 2LTRc with respect to the
 637 mutant D116A. Therefore, the formula applied to estimate the expected levels of 2LTRc with respect to
 638 D116A for sample n is given by $1 - (0.8 \times a_n)$ where a_n is the level of integration measured for sample n,
 639 relative to wt IN A (data from Figure 7A). The values of 2LTRc measured experimentally (Table 1, line
 640 3) are then divided by the expected ones to obtain an estimate of the relative efficiency of nuclear import
 641 (Table 1, line 4).

642 The efficiency of 3' processing was calculated as follows. The ratio of perfect junctions out of
 643 the total amount of 2LTRc (PJ/2LTRc) found for D116A was considered to be the maximal one
 644 and was therefore assigned a value of 1 (Table 1, line 5). The proportion of PJ/2LTRc found for
 645 wt IN A (Table 1) was 0.54 that of D116A (Table 1, line 5). The proportion by which the pool of
 646 PJ/2LTRc found with a catalytically inactive IN can be decreased by 3' processing carried out
 647 by a fully catalytic active IN is therefore $1 - 0.54 = 0.46$ (Table 1, line 6). For mutant NQNK, for
 648 example, the ratio PJ/2LTRc observed was 0.87 of D116A, which corresponds to a relative
 649 decrease of the PJ/2LTRc pool by 0.13 (Table 1, line 6). This decrease is $0.13/0.46 = 0.28$ that
 650 observed for wt IN A (Table 1, line 7), providing an estimate of the relative efficiency of 3'
 651 processing by this mutant with respect to wt IN A. The general formula we applied to estimate

652 the efficiency of 3' processing was therefore $(1-r_x)/0.46$, where 1 is the proportion of PJ/2LTRc
 653 found for D116A, r_x is the ratio PJ/2LTRc observed for sample X and 0.46 is the decrease in
 654 PJ/2LTRc observed for wt IN A with respect to D116A. The resulting values are reported in Table
 655 1, line 7

656 In Table 1 in grey are given the values for the standard deviations (SD). For lines 1, 3 and 5 SD
 657 values are derived from the experimental values; for line 2, SD is given by 1-0.2 multiplied by
 658 the corresponding SD value from line 1; in line 4 SD is given by $((SD_{line3}/average_{line3})^2 +$
 659 $(SD_{line2}/average_{line2})^2)^{1/2} \times average_{line4}$; in line 6 $SD=(SD_{line5}/average_{line5}) \times average_{line6}$; for line 7
 660 $SD=SD_{line6}/0.46$; in line 8 SD is given by $((SD_{line7}/average_{line7})^2 + (SD_{line4}/average_{line4})^2)^{1/2} \times$
 661 $average_{line8}$.

662

663

664 **Assessment of the infectivity of replication-competent viruses**

665 As described above, the coding sequence of NL4.3 IN CTD was replaced with those of wt and mutants
 666 INA CTD. Replication-competent viruses were produced as described above and equal amounts of
 667 viruses were used to infect cells (TZM-bL or CEM-SS) for each sample. For estimating viral replication
 668 in TZM-bL cells, 25 μ L of virus dilution were added to 10^4 cells, plated in 96 wells plates in 75 μ L of
 669 culture medium. After 48h, virus replication was detected by measuring Luc reporter gene expression
 670 by removing 50 μ L of culture medium from each well and adding 50 μ L of Bright Glo reagent to the cells.
 671 After 2 min of incubation at room temperature to allow cell lysis, 100 μ L of cell lysate were transferred to
 672 96-well black solid plates for measurements of luminescence (RLU) using a luminometer (93). For the
 673 detection of virus replication in CEM-SS cells, 0.5×10^6 CEM-SS cells/5ml were infected with 1/25 virus
 674 dilution. After 5 days of culture, the percentage of infected cells were detected by intracellular p24
 675 immuno-staining and flow cytometry analysis as previously described (94).

676

677 **Cloning, production, purification and crystallization**

678 The C-terminal domains (IN CTD, 220-270) of integrases NL4.3, A and A K240Q/N254K studied here
679 were cloned in the pET15b plasmid and the proteins were expressed in BL21DE3 *E. coli* cells. After
680 transformation with the IN C-terminal expressing pET15b, bacteria were inoculated at an OD₆₀₀ of 0.1
681 in one litre of LB medium supplemented with 10% (w/v) sucrose. Cultures were incubated at 37°C with
682 shaking at 220 rpm. At OD_{600nm} of 0.5, the temperature was lowered to 25°C, and shaking reduced to
683 190rpm, till the cells reached an OD_{600nm} of 0.8. IPTG was then added to a final concentration of 0.5 mM
684 to induce the expression of the C-terminal domains. Cells were incubated overnight at 25°. Bacteria
685 were then collected by centrifugation.

686 For protein purification, cells were resuspended in lysis buffer (25 mM HEPES pH 8, 1 M NaCl, 10 mM
687 imidazole) in a ratio of 10 mL of buffer/gram of biomass. Roche Complete Inhibitor Cocktail tablets were
688 added at the beginning of lysis to avoid protease degradation. Cells were lysed by sonification, for
689 1min/g of cells with pulse every 2 seconds at 40% amplitude at 4°C. The bacterial debris were pelleted
690 by ultracentrifugation at 100 000xg for 1hr at 4°C. The supernatant was then loaded on a 1 mL HisTrap
691 FF Crude column (GE Healthcare) with flow rate of 1 mL/min using the AKTA FPLC. Protein was eluted
692 using a gradient up to 500 mM Imidazole in 10 column volumes. Protein concentration was estimated
693 using the Nanodrop. Subsequently, the protein sample was concentrated using the Amicon Ultra 15 mL
694 with a 3 kDa MWCO for the next purification step. A second step of purification was carried out using
695 the S75-16/60 column (GE Healthcare) in 25 mM HEPES pH 8, 1 M NaCl. Samples were dialyzed into
696 25 mM HEPES pH 8, 150 mM NaCl for crystallization.

697 All initial crystallization conditions were determined by vapor diffusion using the TPP Labtech Mosquito
698 Crystal. 200 nL of protein (7-4 mg/mL) was mixed with 200 nL of reservoir in 2 or 3 well of a 96 well
699 MRC crystallization plate which was stored in the Formulatrix RockImager at 20°C. Screen included
700 PEGS (Hampton Research), MPD, CLASSICS, NUCLEIX (Qiagen), JCSG, WIZARDS, ANION and
701 CATION (Molecular Dimensions). Once the initial conditions were obtained, manual drops were set up
702 in Hampton Research 24 well VDX plate to optimize crystallization conditions, and to improve crystal
703 size and quality by mixing 1 µL protein + 1µL reservoir and equilibrating against 500 µL of reservoir at
704 20°C. The IN CTD NL4.3 (group M, subtype B) were obtained in a reservoir containing 0.1 M Tris pH 7,
705 0.8 M potassium sodium tartrate and 0.2 M lithium sulfate. For IN CTD A (subtype A2) and A
706 K240Q/N254K the reservoir was composed of 0.1 M MES pH 6.5 and 1M sodium malonate.

707

708 **Data collection, structure solving and refinement**

709 Data collection was performed at the Swiss Light Source (SLS, Villigen, Switzerland) on a Dectris Pilatus
710 2M detector. After fishing, crystals were rapidly passed through a drop of fluorinated oil (Fomblin® Y
711 LVAC 14/6, average MW 2,500 from Sigma Aldrich) to prevent ice formation and directly frozen on the
712 beamline in the nitrogen stream at 100 K. X-ray diffraction images were indexed and scaled with XDS
713 (95, 96). The structures were solved by molecular replacement using PHASER (97) in the PHENIX (98)
714 program suite using the NMR HIV-1 C-terminal structure (1QMC) (99) as a search model for the IN CTD
715 NL4.3 structure, which was used subsequently as a search model to solve the A and A K240Q/N254K
716 structures. The structure was then built using the AUTOBUILD program (100, 101) followed by several
717 cycles of refinement using PHENIX.REFINE (102) and manual rebuilding with COOT (103). Structure
718 based sequence alignment was performed using PROMALS3D (104). Structures superposition and
719 Root Mean Square Deviations (RMSD) calculations have been performed using secondary structure
720 matching (SSM), superpose program (105) embedded in COOT (103) and in the CCP4 program suite
721 (106). The sequence alignment representation has been generated by ESPript (107). Surface potential
722 was calculated using the DELPHI web server (108) and visualized with CHIMERA (109).
723 Crystallographic structures were deposited in PDB under the identification numbers 6T6E (HIV-1 Cter,
724 PNL4.3), 6T6I (HIV-1 Cter, subtype A2) and 6T6J (HIV-1 Cter, subtype A2, mutant N254K-K240Q).

725

726 **Analysis of the surface electrostatic potential of retroviral CTDs**

727 The structures and the sequences of the C-terminal domains have been extracted from: HIV-1 A2, PDB
728 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V (48); MVV, PDB
729 5LLJ (49); RSV, PDB 1COM (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV, PDB 4E7I
730 (53). The structure based sequence alignment has been performed using PROMALS3D (104).
731 Structures superposition and rRoot Mean Square Deviations (RMSD) calculations have been performed
732 using secondary structure matching (SSM), superpose program (105) embedded in COOT (103). The

733 sequence alignment representation has been generated by ESPript (107). The surface electrostatic
734 potential was calculated using the DELPHI web server (108) and visualized with CHIMERA (109).

735

736 **Statistical tests**

737 All statistical analyses were performed on at least three independent experiments (transfection and
738 transduction) using Prism 6 (GraphPad). For all functional tests, the values obtained for the chimeras
739 were normalized using the values obtained for parental integrase A. Student tests were used to evaluate
740 whether the normalized mean values obtained with the chimeric and mutant integrases were
741 significantly different from that obtained with the parental strain, and/or between them. For confocal
742 microscopy, unpaired t test was used for statistical analyses.

743

744 **Data availability**

745 Crystallographic structures are available in PDB under identification numbers 6T6E, 6T6I, and 6T6J.

746

747 **Acknowledgments**

748 The authors are grateful to Pr. J.C. Plantier for providing HIV-1 strains of subtype B, CRF02 and group
749 O, to C. Elefante for the construction of the p8.91-MB plasmid, to J. Batisse for providing control
750 reagents, and to M. Lavigne, B. Maillot and S. Marzi for helpful discussions. The authors wish to thank
751 R. Drillien (IGBMC) for suggestions about the manuscript. The authors thank V. Olieric and the staff of
752 the Swiss Light Source synchrotron for help with data collection. The authors acknowledge the support
753 and the use of resources of the French Infrastructure for Integrated Structural Biology FRISBI ANR-10-
754 INBS-05 and of Instruct-ERIC.

755

756

757 **Table and figure legends**

758 **Table 1. Estimate of the contribution of the defects of nuclear import and 3' processing to the**
 759 **efficiency of integration observed with the mutants of the NKNK motif.**

760 **Table 2. Primers and probes used in the qPCR assay.**

761 **Table 3. Primers used for the Alu qPCR assay.**

762 **Table 4. Primers and probes used for the detection of total 2LTRc and 2LTRc with perfect**
 763 **palindromic junction (PJ).**

764 **Figure 1. Outline of the experimental system.** *Panel A.* Workflow used to evaluate Pr55Gag
 765 processing, reverse transcription, and integration in our experimental system. VSV-pseudotyped HIV-1
 766 derived vectors, produced by triple transfection, were used to transduce HEK 293T cells. Upon
 767 integration, the proviral DNA will allow growth of the cellular clones in the presence of puromycin. For
 768 multiplicities of infection lower than 1, the number of clones obtained is directly proportional to the
 769 number of integration events. *Panel B.* schematic representation of the viral genomic RNA contained in
 770 the viral vectors, transcribed from pSRP (panel A, also see Materials and Methods). R, U5 and U3, viral
 771 sequences constituting the LTR; "cis-acting", viral sequences required for RNA packaging and reverse
 772 transcription; EF1- α and hPGK, internal human promoters driving the expression of the nuclear RFP
 773 (nRFP) and of the puromycin N-acetyl-transferase that confers resistance to puromycin (Puro^R),
 774 respectively. *Panel C.* Evaluation of reverse transcription and integration in control samples, compared
 775 to v8.91-MB reference vector. The results give the average values of three independent experiments.

776 **Figure 2: Functionality of chimerical integrases.** *Panel A.* Alignment of CTD sequences from isolates
 777 A and O, used in this study. The numbers in italic on the left and on the right of the alignment indicate
 778 the beginning and the end (in amino acid) of the CTD, respectively. Only amino acids divergent between
 779 the two sequences are indicated by letters. The arrows and numbers above the alignment indicate the
 780 last position that, in the chimeras, was concordant with the sequence of isolate A. *Panel B.* Schematic
 781 representation of the integrases studied. Integrase from isolate O is drawn at the top of the panel in dark
 782 grey; integrase from isolate A is drawn at the bottom of the panel in light grey. The genetic origin of the

783 portions of the chimeras is indicated by the colour code that refers to the reference isolates A and O.
 784 *Panel C.* Representative western blot obtained with an anti-CA mouse monoclonal antibody. *Panel D.*
 785 Efficiency of processing of the Pr55Gag precursor, estimated by the amount of CA compared to the
 786 amount of Pr55Gag precursors detected by western blot (as in panel B). The results are expressed as
 787 function of the reference wt IN A, set at 100%. *Panel E.* Efficiency of reverse transcription (detection of
 788 the junction U5-Psi by qPCR) expressed as function of the reference wt IN A. *Panel F.* Efficiency of
 789 integration calculated with the puromycin assay, normalized by the amount of total viral DNA (estimated
 790 by qPCR), expressed as function of the reference wt IN A. Error bars indicate standard deviations. The
 791 results given in panels C-E are the average of 3 independent experiments. ** p <0.01; *** p <0.001, p
 792 values for comparison to wt IN A.

793 **Figure 3. Functionality of IN A with mutated CTD.** *Panels A-C.* Efficiency of processing of the
 794 Pr55Gag precursor (panel A), of reverse transcription (panel B), and normalized efficiency of integration
 795 (panel C). *Panel D.* Efficiency of processing of the Pr55Gag precursor, of reverse transcription, and of
 796 normalized efficiency of integration for the K240Q/K273Q mutant. *Panel E.* Efficiency of integration of
 797 the K240R/K273R mutant (NRNR in the Figure) and of wt IN A (*NKNK, set at 100%). Error bars indicate
 798 standard deviations. In all panels the results are the average of 3 independent experiments. * p <0.05;
 799 ** p <0.01; *** p <0.001, p values for comparison to wt IN A.

800 **Figure 4. Definition of the NKNK motif and of its importance in the most widespread phylogenetic**
 801 **groups of HIV-1.** *Panel A.* Efficiency of integration of various mutants of the N residues of the NKNK
 802 motif and of the wt enzyme (*NKNK, set at 100%). *Panel B.* Efficiency of integration of the mutant
 803 carrying the sequences of isolate O at positions 222, 240, 254 and 273 (K₂₂₂Q₂₄₀K₂₅₄Q₂₇₃, KQKQ in the
 804 Figure) and of the wt enzyme (*NKNK, set at 100%). *Panel C.* Conservation logo of the sequence at
 805 positions 222, 240, 254 and 273 in HIV-1 group M integrases. *Panel D.* Efficiency of integration of the
 806 double mutant N240Q/K273Q (NQNQ in the Figure) of an isolate of subtype B, one of subtype C and
 807 from CRF02, compared to the corresponding wt integrases, set as reference at 100%. In all vectors the
 808 RT sequence had the same phylogenetic origin as IN and was replaced using the *MluI-BspEI* cassette
 809 in p8.91MB, as described in Materials and Methods. Error bars indicate standard deviations (standard
 810 deviations of the wt integrases of each subtype were calculated with respect to reference wt IN A, used
 811 as control). The results are the average of 3 independent experiments.

812 **Figure 5. Importance of the number and position of the K residues in the N₂₂₂K₂₄₀N₂₅₄K₂₇₃ motif**
 813 **of the CTD.** *Panel A.* Efficiency of integration, normalized by the amount of viral DNA, for the IN mutants
 814 grouped by the number of K present at positions 222, 240, 254, 273. The composition in amino acids in
 815 the four positions of the motif is given for the isolate with 0 and for the one with 4 K. For clarity, only the
 816 four letters of the amino acids of the motif are represented for each mutant, omitting the positions; the
 817 first letter indicates the residue at position 222, the second, position 240, the third, position 254 and the
 818 fourth, position 273. *Panels B-D.* Efficiency of integration of the individual mutants containing 1 (panel
 819 B), 2 (panel C), or 3 (panel D) K in the motif. In panel C, the motif corresponding to the sequence of wt
 820 IN A (reference set at 100%) is indicated by an asterisk. Error bars indicate standard deviations. The
 821 results are the average of 4 independent experiments. * p <0.05; ** p <0.01; *** p <0.001, p values for
 822 comparison to wt IN A. *Panel E.* Efficiency of reverse transcription for the mutants containing one, two,
 823 three, or four K in the motif. The motif of wt IN A is indicated by an asterisk. The color of the bars reflects
 824 the colors used in figure 5, panels A, B, C and D. The results are the average of at least 3 independent
 825 experiments: * p <0.05; ** p <0.01; *** p <0.001, p values for comparison to wt IN A.

826 **Figure 6. Importance of the NKNK motif in replication-competent viruses.** *Panel A.* Scheme of the
 827 portion coding for the integrase in the various viruses. Drawn in grey are the parts derived from the
 828 NL4.3 sequences, in white those from isolate A. The black bars indicate positions 222, 240, 254 and
 829 273 from left to right; the amino acid found for each mutant in each of these four positions is indicated
 830 above the bars. *Panel B.* Infectivity of the viruses shown in panel E (except wt NL4.3 that is used as
 831 reference, set at 100%). The results are given in grey for CEM-SS and in black for TZM-bL cells. Error
 832 bars indicate standard deviations with respect to the reference wt pNL4-3. The results are the average
 833 of 2 independent experiments.

834 **Figure 7. Amount of 2LTRc/total viral DNA (panel A) and of ratio of PJ/2LTRc (panel B) in the**
 835 **mutants deprived of one or both K of the NKNK motif.** The motif corresponding to the sequence of
 836 wt IN A (reference set at 100%) is indicated by an asterisk. Error bars indicate standard deviations.
 837 Above the plot are given the p values for the comparisons of the different samples with respect to wt IN
 838 A or to the integration-deficient mutant IN A D116A (* p <0.05; ** p <0.01; *** p <0.001). The number of
 839 independent experiments performed for each sample (n) is also given. The amount of 2LTRc for sample
 840 X is normalised by the total amount of viral DNA (detection of the late reverse transcription products by

841 quantifying the U5-Psi junction), then it is expressed as a percentage of the amount detected for the
 842 control INA D116A (indicated with D), thus giving $(2LTR_{Cx}/r_x)/(2LTR_{CD}/r_D) \times 100$, where r_x and r_D are the
 843 amounts of late reverse transcription products, in sample X and in control D, respectively, and $2LTR_{Cx}$
 844 and $2LTR_{CD}$ the amount of 2LTR circles in sample X and in control D, respectively.

845 **Figure 8. Structural analysis of the NKNK motif.** *Panels A and B.* Side view (A) and top view (B) of
 846 the ribbon representation of the crystal structure of CTD A. The positions of residues N222, K240 and
 847 N254 are represented with sticks ~~as well as the position of the I234, the only different residue between~~
 848 ~~the CTD of IN A and IN NL4.3.~~ *Panels C to F.* Surface electrostatic potential representation of the CTD
 849 A wt (panel C side view and panel D top view) and of the NQKK mutant (panel E side view and panel F
 850 top view). In red, negative potential; in blue, positive potential and in white neutral regions. Circled in
 851 yellow is the region with large differences in the mutant structures ~~(see below Panel H-M).~~ *Panel G.*
 852 Superposition of the CTDs of IN A and IN A NQKK (chains A, B and C). The mutation N254K induces a
 853 displacement of the K236 side chain (white arrows) disturbing the structure of the 235-237 region.
 854 ~~*Panels F and G.* Side view (F) and top view (G) of the superposition of the three molecules in the~~
 855 ~~asymmetric unit of the NQKK CTD. The position of the residue N222, K240Q mutation and N254K~~
 856 ~~mutation are represented with sticks as well as the position of the I234.~~ *Panels H, J and L (side view)*
 857 ~~and I, K and M (top view) are the surface electrostatic potential representation of NQKK CTD chain A~~
 858 ~~(H, I), chain B (J, K) and chain C (L, M). In red, negative potential; in blue, positive potential and in white~~
 859 ~~neutral regions. The regions with large differences are circled in yellow.~~

860 **Figure 9. Analysis of the surface electrostatic potential in the C-terminal of retroviral integrases.**
 861 The structures and the sequences of the C-terminal domains have been extracted from: HIV-1 A2, PDB
 862 6T6I (this publication); HIV-1 PNL4.3, PDB 6T6E (this publication); SIV, PDB 1C6V (48); MVV, PDB
 863 5LLJ (49); RSV, PDB 1COM (50); MMTV, PDB 5D7U (51); MMLV, PDB 2M9U (52); PFV, PDB 4E7I
 864 (53). *Panel A.* Superposition of the structure of the integrase C-terminal domains of four lentiviruses
 865 (HIV-1 A2, pink; HIV-1 pNL4.3, orange; Simian Immunodeficiency Virus, SIV, khaki; Maedi-Visna virus,
 866 MVV, cyan), of an α retrovirus (the Rous Sarcoma Virus, RSV, blue), of a β retrovirus (the Mouse
 867 Mammary Tumor Virus, MMTV, sky blue), of a γ retrovirus (the Moloney Murine Leukemia Virus, MMLV,
 868 purple), and of a spumaretrovirus (the Prototype Foamy Virus, PFV, green). *Panel B.* Structure-based
 869 sequences alignment of the integrases C-terminal domains. Sequence numbering corresponds to the

870 HIV-1 A2 integrase sequence. Secondary structures from HIV-1 A2 are represented (TT: β -Turn, β 1 to
871 β 5: β -sheets, η 1: 3_{10} -helix). Residues framed in blue: Position in the alignment of the three first amino
872 acids from the NKNK motif. Red background: 100% identity in the sequence alignment. Yellow
873 background: % of equivalent residues > 70% (considering their physical-chemical properties),
874 equivalent residues are depicted in bold. *Panels C-G*. Surface electrostatic potential representation of
875 integrases from several retroviruses. The surface corresponding to that delimited by the NKNK motif
876 in HIV-1 M (~~panels C, D and E~~) is circled in yellow. *Panel CO*, ribbon representation of HIV-1 A2 C-
877 terminal domain structure. The amino acids belonging to the NKNK motif are represented in sticks.
878 *Panels C1-C4*. Surface potential representation of the C-terminal domain structures of four lentiviral
879 integrases: HIV-1 A2 (panel C1), HIV-1 pNL4.3 (panel C2), Simian Immunodeficiency Virus (SIV) (panel
880 C3) and Maedi-Visna virus (MVV) (panel C4). *Panel D0*. Ribbon representation of the structure of Rous
881 Sarcoma Virus (RSV) C-terminal domain. The surface circled in yellow corresponds to that delimited
882 by the NKNK motif in HIV-1 M after superposition of the structures. The amino acids corresponding to
883 the motif in the structure-based alignment are shown as sticks. *Panel D1*. Surface potential
884 representation of the C-terminal domain structure of an α retrovirus, the Rous Sarcoma Virus (RSV).
885 *Panel E*. Surface potential representation of the C-terminal domain structure of a β retrovirus, the Mouse
886 Mammary Tumor Virus (MMTV). *Panel F*. Surface potential representation of the C-terminal domain
887 structure of a γ retrovirus, the Moloney Murine Leukemia Virus (MMLV). *Panel G*. Surface potential
888 representation of the C-terminal domain structure of a spumaretrovirus, the Prototype Foamy Virus
889 (PFV). Negative potential is in red, neutral in white and positive potential is in blue.

890 **Figure 10. The NKNK motif in the post-integration intasome structure.** Two 90° views of the
891 intasome structure from Passos et al. (31). The residues forming the electrostatic surface are pointing
892 towards the solvent for the two independant C-terminal domains of the integrase dimer

893

894

895

896

897

898

References

- 899 1. **Pauza CD.** 1990. Two bases are deleted from the termini of HIV-1 linear DNA during
900 integrative recombination. *Virology* **179**:886–889.
- 901 2. **Engelman A, Mizuuchi K, Craigie R.** 1991. HIV-1 DNA Integration - Mechanism of Viral-
902 Dna Cleavage and Dna Strand Transfer. *Cell* **67**:1211–1221.
- 903 3. **Engelman A, Englund G, Orenstein JM, Martin MA, Craigie R.** 1995. Multiple Effects of
904 Mutations in Human-Immunodeficiency-Virus Type-1 Integrase on Viral Replication. *J Virol*
905 **69**:2729–2736.
- 906 4. **Bukovsky A, Gottlinger H.** 1996. Lack of integrase can markedly affect human
907 immunodeficiency virus type 1 particle production in the presence of an active viral protease.
908 *J Virol* **70**:6820–6825.
- 909 5. **Hoyte AC, Jamin AV, Koneru PC, Kobe MJ, Larue RC, Fuchs JR, Engelman AN,**
910 **Kvaratskhelia M.** 2017. Resistance to pyridine-based inhibitor KF116 reveals an unexpected
911 role of integrase in HIV-1 Gag-Pol polyprotein proteolytic processing. *J Biol Chem*
912 **292**:19814–19825.
- 913 6. **Kessi JJ, Kutluay SB, Townsend D, Rebensburg S, Slaughter A, Larue RC, Shkriabai N,**
914 **Bakouche N, Fuchs JR, Bieniasz PD, Kvaratskhelia M.** 2016. HIV-1 Integrase binds the
915 viral RNA genome and is essential during virion morphogenesis. *Cell* **166**:1257–1268.e12.
- 916 7. **Zhu K, Dobard C, Chow SA.** 2004. Requirement for integrase during reverse transcription of
917 human immunodeficiency virus type 1 and the effect of cysteine mutations of integrase on its
918 interactions with reverse transcriptase. *J Virol* **78**:5045–5055.
- 919 8. **Dobard CW, Briones MS, Chow SA.** 2007. Molecular mechanisms by which human
920 immunodeficiency virus type 1 integrase stimulates the early steps of reverse transcription. *J*
921 *Virol* **81**:10037–10046.
- 922 9. **Wilkinson TA, Januszyk K, Phillips ML, Tekeste SS, Zhang M, Miller JT, Le Grice SFJ,**
923 **Clubb RT, Chow SA.** 2009. Identifying and characterizing a functional HIV-1 reverse
924 transcriptase-binding site on integrase. *J Biol Chem* **284**:7931–7939.
- 925 10. **Emiliani S, Mousnier A, Busschots K, Maroun M, Van Maele B, Tempé D,**
926 **Vandekerckhove L, Moisant F, Ben-Slama L, Witvrouw M, Christ F, Rain J-C,**
927 **Dargemont C, Debyser Z, Benarous R.** 2005. Integrase mutants defective for interaction
928 with LEDGF/p75 are impaired in chromosome tethering and HIV-1 replication. *J Biol Chem*
929 **280**:25517–25523.
- 930 11. **Bukrinsky MI, Sharova N, Dempsey MP, Stanwick TL, Bukrinskaya AG, Haggerty S,**
931 **Stevenson M.** 1992. Active nuclear import of human immunodeficiency virus type 1
932 preintegration complexes. *PNAS* **89**:6580–6584.
- 933 12. **Mattaj JW, Englmeier L.** 1998. Nucleocytoplasmic transport: the soluble phase. *Annu Rev*
934 *Biochem* **67**:265–306.
- 935 13. **Yamashita M, Emerman M.** 2004. Capsid is a dominant determinant of retrovirus infectivity
936 in nondividing cells. *J Virol* **78**:5670–5678.
- 937 14. **Yamashita M, Perez O, Hope TJ, Emerman M.** 2007. Evidence for Direct Involvement of
938 the Capsid Protein in HIV Infection of Nondividing Cells. *PLoS Pathog* **3**:e156.
- 939 15. **Di Nunzio F, Danckaert A, Fricke T, Perez P, Fernandez J, Perret E, Roux P, Shorte S,**
940 **Charneau P, Diaz-Griffero F, Arhel NJ.** 2012. Human nucleoporins promote HIV-1 docking
941 at the nuclear pore, nuclear import and integration. *PLoS ONE* **7**:e46037.

- 942 16. **Di Nunzio F, Fricke T, Miccio A, Valle-Casuso JC, Perez P, Souque P, Rizzi E,**
943 **Severgnini M, Mavilio F, Charneau P, Diaz-Griffero F.** 2013. Nup153 and Nup98 bind the
944 HIV-1 core and contribute to the early steps of HIV-1 replication. *Virology* **440**:8–18.
- 945 17. **Matreyek KA, Engelman A.** 2011. The requirement for nucleoporin NUP153 during human
946 immunodeficiency virus type 1 infection is determined by the viral capsid. *J Virol* **85**:7818–
947 7827.
- 948 18. **Krishnan L, Matreyek KA, Oztop I, Lee K, Tipper CH, Li X, Dar MJ, KewalRamani VN,**
949 **Engelman A.** 2010. The requirement for cellular transportin 3 (TNPO3 or TRN-SR2) during
950 infection maps to human immunodeficiency virus type 1 capsid and not integrase. *J Virol*
951 **84**:397–406.
- 952 19. **Cribrier A, Ségéral E, Delelis O, Parissi V, Simon A, Ruff M, Benarous R, Emiliani S.**
953 2011. Mutations affecting interaction of integrase with TNPO3 do not prevent HIV-1 cDNA
954 nuclear import. *Retrovirology* **8**:104.
- 955 20. **Gallay P, Hope T, Chin D, Trono D.** 1997. HIV-1 infection of nondividing cells through the
956 recognition of integrase by the importin/karyopherin pathway. *PNAS* **94**:9825–9830.
- 957 21. **Bouyac-Bertoia M, Dvorin JD, Fouchier RAM, Jenkins Y, Meyer BE, Wu LI, Emerman M,**
958 **Malim MH.** 2001. HIV-1 infection requires a functional integrase NLS. *Molecular Cell* **7**:1025–
959 1035.
- 960 22. **Hearps AC, Jans DA.** 2006. HIV-1 integrase is capable of targeting DNA to the nucleus via
961 an Importin α/β -dependent mechanism. *Biochemical Journal* **398**:475–484.
- 962 23. **Ao Z, Huang G, Yao H, Xu Z, Labine M, Cochrane AW, Yao X.** 2007. Interaction of human
963 immunodeficiency virus type 1 integrase with cellular nuclear import receptor importin 7 and
964 its impact on viral replication. *J Biol Chem* **282**:13456–13467.
- 965 24. **Ao Z, Jayappa KD, Wang B, Zheng Y, Kung S, Rassart E, Depping R, Kohler M, Cohen**
966 **EA, Yao X.** 2010. Importin $\alpha 3$ Interacts with HIV-1 Integrase and Contributes to HIV-1
967 Nuclear Import and Replication. *J Virol* **84**:8650–8663.
- 968 25. **Woodward CL, Prakobwanakit S, Mosessian S, Chow SA.** 2009. Integrase interacts with
969 nucleoporin NUP153 to mediate the nuclear import of human immunodeficiency virus type 1.
970 *J Virol* **83**:6522–6533.
- 971 26. **Ao Z, Jayappa KD, Wang B, Zheng Y, Wang X, Peng J, Yao X.** 2012. Contribution of host
972 nucleoporin 62 in HIV-1 integrase chromatin association and viral DNA integration. *J Biol*
973 *Chem* **287**:10544–10555.
- 974 27. **Larue R, Gupta K, Wuensch C, Shkriabai N, Kessl JJ, Danhart E, Feng L, Taltynov O,**
975 **Christ F, Van Duynne GD, Debyser Z, Foster MP, Kvaratskhelia M.** 2012. Interaction of the
976 HIV-1 intasome with transportin 3 protein (TNPO3 or TRN-SR2). *J Biol Chem* **287**:34044–
977 34058.
- 978 28. **De Houwer S, Demeulemeester J, Thys W, Rocha S, Dirix L, Gijssbers R, Christ F,**
979 **Debyser Z.** 2014. The HIV-1 integrase mutant R263A/K264A is 2-fold defective for TRN-SR2
980 binding and viral nuclear import. *J Biol Chem* **289**:25351–25361.
- 981 29. **Christ F, Thys W, De Rijck J, Gijssbers R, Albanese A, Arosio D, Emiliani S, Rain J-C,**
982 **Benarous R, Cereseto A, Debyser Z.** 2008. Transportin-SR2 imports HIV into the nucleus.
983 *Curr Biol* **18**:1192–1202.
- 984 30. **Jayappa KD, Ao Z, Yang M, Wang J, Yao X.** 2011. Identification of critical motifs within HIV-
985 1 integrase required for importin $\alpha 3$ interaction and viral cDNA nuclear import. *J Mol Biol*
986 **410**:847–862.

- 987 31. **Passos DO, Li M, Yang R, Rebensburg SV, Ghirlando R, Jeon Y, Shkriabai N,**
988 **Kvaratskhelia M, Craigie R, Lyumkis D.** 2017. Cryo-EM structures and atomic model of the
989 HIV-1 strand transfer complex intasome. *Science* **355**:89–92.
- 990 32. **Michel F, Crucifix C, Granger F, Eiler S, Mouscadet J-F, Korolev S, Agapkina J,**
991 **Ziganshin R, Gottikh M, Nazabal A, Emiliani S, Benarous R, Moras D, Schultz P, Ruff M.**
992 2009. Structural basis for HIV-1 DNA integration in the human genome, role of the
993 LEDGF/P75 cofactor. *EMBO J* **28**:980–991.
- 994 33. **Craigie R, Bushman FD.** 2012. HIV DNA integration. *Cold Spring Harb Perspect Med*
995 **2**:a006890–a006890.
- 996 34. **Delelis O, Carayon K, Saïb A, Deprez E, Mouscadet J-F.** 2008. Integrase and integration:
997 biochemical activities of HIV-1 integrase. *Retrovirology* **5**:114.
- 998 35. **Zheng R, Jenkins TM, Craigie R.** 1996. Zinc folds the N-terminal domain of HIV-1 integrase,
999 promotes multimerization, and enhances catalytic activity. *Proc Natl Acad Sci USA*
1000 **93**:13659–13664.
- 1001 36. **Eijkelenboom AP, van den Ent FM, Vos A, Doreleijers JF, Hård K, Tullius TD, Plasterk**
1002 **RH, Kaptein R, Boelens R.** 1997. The solution structure of the amino-terminal HHCC
1003 domain of HIV-2 integrase: a three-helix bundle stabilized by zinc. *Curr Biol* **7**:739–746.
- 1004 37. **Busschots K, Vercammen J, Emiliani S, Benarous R, Engelborghs Y, Christ F, Debyser**
1005 **Z.** 2005. The interaction of LEDGF/p75 with integrase is lentivirus-specific and promotes
1006 DNA binding. *J Biol Chem* **280**:17841–17847.
- 1007 38. **Heuer TS, Brown PO.** 1997. Mapping features of HIV-1 integrase near selected sites on viral
1008 and target DNA molecules in an active enzyme-DNA complex by photo-cross-linking.
1009 *Biochemistry* **36**:10655–10665.
- 1010 39. **Esposito D, Craigie R.** 1998. Sequence specificity of viral end DNA binding by HIV-1
1011 integrase reveals critical regions for protein-DNA interaction. *EMBO J* **17**:5832–5843.
- 1012 40. **Chen A, Weber IT, Harrison RW, Leis J.** 2006. Identification of amino acids in HIV-1 and
1013 avian sarcoma virus integrase subsites required for specific recognition of the long terminal
1014 repeat Ends. *J Biol Chem* **281**:4173–4182.
- 1015 41. **Engelman A, Hickman AB, Craigie R.** 1994. The core and carboxyl-terminal domains of the
1016 integrase protein of human immunodeficiency virus type 1 each contribute to nonspecific
1017 DNA binding. *J Virol* **68**:5911–5917.
- 1018 42. **Lutzke RA, Vink C, Plasterk RH.** 1994. Characterization of the minimal DNA-binding
1019 domain of the HIV integrase protein. *Nucleic Acids Res* **22**:4125–4131.
- 1020 43. **Cannon PM, Byles ED, Kingsman SM, Kingsman AJ.** 1996. Conserved sequences in the
1021 carboxyl terminus of integrase that are essential for human immunodeficiency virus type 1
1022 replication. *J Virol* **70**:651–657.
- 1023 44. **Larder BA, Purifoy DJ, Powell KL, Darby G.** 1987. Site-specific mutagenesis of AIDS virus
1024 reverse transcriptase. *Nature* **327**:716–717.
- 1025 45. **Mandal D, Feng Z, Stoltzfus CM.** 2008. Gag-processing defect of human immunodeficiency
1026 virus type 1 integrase E246 and G247 mutants is caused by activation of an overlapping 5'
1027 splice site. *J Virol* **82**:1600–1604.

- 1028 46. **Hemelaar J, Gouws E, Ghys PD, Osmanov S, WHO-UNAIDS Network for HIV Isolation and Characterisation.** 2011. Global trends in molecular epidemiology of HIV-1 during 2000-2007. *AIDS* **25**:679–689.
- 1031 47. **Sloan RD, Wainberg MA.** 2011. The role of unintegrated DNA in HIV infection. *Retrovirology* **8**:52.
- 1033 48. **Chen Z, Yan Y, Munshi S, Li Y, Zugay-Murphy J, Xu B, Witmer M, Felock P, Wolfe A, Sardana V, Emini EA, Hazuda D, Kuo LC.** 2000. X-ray structure of simian immunodeficiency virus integrase containing the core and C-terminal domain (residues 50-293): an initial glance of the viral DNA binding platform. *J Mol Biol* **296**:521–533.
- 1037 49. **Ballandras-Colas A, Maskell DP, Serrao E, Locke J, Swuec P, Jónsson SR, Kotecha A, Cook NJ, Pye VE, Taylor IA, Andrésdóttir V, Engelman AN, Costa A, Cherepanov P.** 2017. A supramolecular assembly mediates lentiviral DNA integration. *Science* **355**:93–95.
- 1040 50. **Yang ZN, Mueser TC, Bushman FD, Hyde CC.** 2000. Crystal structure of an active two-domain derivative of Rous sarcoma virus integrase. *J Mol Biol* **296**:535–548.
- 1042 51. **Ballandras-Colas A, Brown M, Cook NJ, Dewdney TG, Demeler B, Cherepanov P, Lyumkis D, Engelman AN.** 2016. Cryo-EM reveals a novel octameric integrase structure for betaretroviral intasome function. *Nature* **530**:358–361.
- 1045 52. **Aiyer S, Swapna GVT, Malani N, Aramini JM, Schneider WM, Plumb MR, Ghanem M, Larue RC, Sharma A, Studamire B, Kvaratskhelia M, Bushman FD, Montelione GT, Roth MJ.** 2014. Altering murine leukemia virus integration through disruption of the integrase and BET protein family interaction. *Nucleic Acids Res* **42**:5917–5928.
- 1049 53. **Hare S, Maertens GN, Cherepanov P.** 2012. 3'-processing and strand transfer catalysed by retroviral integrase in crystallo. *EMBO J* **31**:3020–3028.
- 1051 54. **Luscombe NM, Laskowski RA, Thornton JM.** 2001. Amino acid-base interactions: a three-dimensional analysis of protein-DNA interactions at an atomic level. *Nucleic Acids Res* **29**:2860–2874.
- 1054 55. **Zhou H-X, Pang X.** 2018. Electrostatic interactions in protein structure, folding, binding, and condensation. *Chem Rev* **118**:1691–1741.
- 1056 56. **Karlin S, Zhu ZY.** 1996. Characterizations of diverse residue clusters in protein three-dimensional structures. *PNAS* **93**:8344–8349.
- 1058 57. **Karlin S, Brendel V.** 1988. Charge configurations in viral proteins. *PNAS* **85**:9396–9400.
- 1059 58. **Parker MS, Balasubramaniam A, Parker SL.** 2012. On the segregation of protein ionic residues by charge type. *Amino Acids* **43**:2231–2247.
- 1061 59. **Kharrat N, Belmabrouk S, Abdelhedi R, Benmarzoug R, Assidi M, Qahtani AI MH, Rebai A.** 2016. Screening for clusters of charge in human virus proteomes. *BMC Genomics* **17**:758–19.
- 1064 60. **Dostálková A, Kaufman F, Křížová I, Kultová A, Strohalmová K, Hadravová R, Ruml T, Rumlová M.** 2018. Mutations in the basic region of the Mason-Pfizer Monkey virus nucleocapsid protein affect reverse transcription, genomic RNA packaging, and the virus assembly site. *J Virol* **92**:5439.
- 1068 61. **Hamann MV, Müllers E, Reh J, Stanke N, Effantin G, Weissenhorn W, Lindemann D.** 2014. The cooperative function of arginine residues in the Prototype Foamy Virus Gag C-terminus mediates viral and cellular RNA encapsidation. *Retrovirology*, 6 ed. **11**:87–17.

- 1071 62. **Heyrana KJ, Goh BC, Perilla JR, Nguyen T-LN, England MR, Bewley MC, Schulten K,**
1072 **Craven RC.** 2016. Contributions of charged residues in structurally dynamic capsid surface
1073 loops to Rous Sarcoma virus assembly. *J Virol* **90**:5700–5714.
- 1074 63. **Yenamandra SP, Sompallae R, Klein G, Kashuba E.** 2009. Comparative analysis of the
1075 Epstein-Barr virus encoded nuclear proteins of EBNA-3 family. *Comput Biol Med* **39**:1036–
1076 1042.
- 1077 64. **Roller RJ, Bjerke SL, Haugo AC, Hanson S.** 2010. Analysis of a charge cluster mutation of
1078 herpes simplex virus type 1 UL34 and its extragenic suppressor suggests a novel interaction
1079 between pUL34 and pUL31 that is necessary for membrane curvature around capsids. *J Virol*
1080 **84**:3921–3934.
- 1081 65. **Zhang J, Webb DM.** 2004. Rapid evolution of primate antiviral enzyme APOBEC3G. *Hum*
1082 *Mol Genet* **13**:1785–1791.
- 1083 66. **Vasudev PG, Banerjee M, Ramakrishnan C, Balaram P.** 2012. Asparagine and glutamine
1084 differ in their propensities to form specific side chain-backbone hydrogen bonded motifs in
1085 proteins. *Proteins* **80**:991–1002.
- 1086 67. **Weichenberger CX, Sippl MJ.** 2006. Self-consistent assignment of asparagine and
1087 glutamine amide rotamers in protein crystal structures. *Structure* **14**:967–972.
- 1088 68. **Métifiot M, Johnson BC, Kiselev E, Marler L, Zhao XZ, Burke TR, Marchand C, Hughes**
1089 **SH, Pommier Y.** 2016. Selectivity for strand-transfer over 3'-processing and susceptibility to
1090 clinical resistance of HIV-1 integrase inhibitors are driven by key enzyme-DNA interactions in
1091 the active site. *Nucleic Acids Res* **44**:6896–6906.
- 1092 69. **Johnson AA, Santos W, Pais GCG, Marchand C, Amin R, Burke TR, Verdine G,**
1093 **Pommier Y.** 2006. Integration requires a specific interaction of the donor DNA terminal 5'-
1094 cytosine with glutamine 148 of the HIV-1 integrase flexible loop. *J Biol Chem* **281**:461–467.
- 1095 70. **Guiot E, Carayon K, Delelis O, Simon F, Tauc P, Zubin E, Gottikh M, Mouscadet J-F,**
1096 **Brochon J-C, Deprez E.** 2006. Relationship between the oligomeric status of HIV-1
1097 integrase on DNA and enzymatic activity. *J Biol Chem* **281**:22707–22719.
- 1098 71. **Hehl EA, Joshi P, Kalpana GV, Prasad VR.** 2004. Interaction between human
1099 immunodeficiency virus type 1 reverse transcriptase and integrase proteins. *J Virol* **78**:5056–
1100 5067.
- 1101 72. **Ao Z, Fowke KR, Cohen EA, Yao X.** 2005. Contribution of the C-terminal tri-lysine regions
1102 of human immunodeficiency virus type 1 integrase for efficient reverse transcription and viral
1103 DNA nuclear import. *Retrovirology* **2**:62.
- 1104 73. **Williams KL, Zhang Y, Shkriabai N, Karki RG, Nicklaus MC, Kotrikadze N, Hess S, Le**
1105 **Grice SFJ, Craigie R, Pathak VK, Kvaratskhelia M.** 2005. Mass spectrometric analysis of
1106 the HIV-1 integrase-pyridoxal 5'-phosphate complex reveals a new binding site for a
1107 nucleotide inhibitor. *J Biol Chem* **280**:7949–7955.
- 1108 74. **Tekeste SS, Wilkinson TA, Weiner EM, Xu X, Miller JT, Le Grice SFJ, Clubb RT, Chow**
1109 **SA.** 2015. Interaction between Reverse Transcriptase and Integrase Is Required for Reverse
1110 Transcription during HIV-1 Replication. *J Virol* **89**:12058–12069.
- 1111 75. **Dar MJ, Monel B, Krishnan L, Shun M-C, Di Nunzio F, Helland DE, Engelman A.** 2009.
1112 Biochemical and virological analysis of the 18-residue C-terminal tail of HIV-1 integrase.
1113 *Retrovirology* **6**:94.

- 1114 76. **Mohammed KD, Topper MB, Muesing MA.** 2011. Sequential deletion of the integrase
1115 (Gag-Pol) carboxyl terminus reveals distinct phenotypic classes of defective HIV-1. *J Virol*
1116 **85**:4654–4666.
- 1117 77. **Lu R, Ghory HZ, Engelman A.** 2005. Genetic analyses of conserved residues in the
1118 carboxyl-terminal domain of human immunodeficiency virus type 1 integrase. *J Virol*
1119 **79**:10356–10368.
- 1120 78. **Cereseto A, Manganaro L, Gutierrez MI, Terreni M, Fittipaldi A, Lusic M, Marcello A,**
1121 **Giacca M.** 2005. Acetylation of HIV-1 integrase by p300 regulates viral integration. *EMBO J*
1122 **24**:3070–3081.
- 1123 79. **Terreni M, Valentini P, Liverani V, Gutierrez MI, Di Primio C, Di Fenza A, Tozzini V,**
1124 **Allouch A, Albanese A, Giacca M, Cereseto A.** 2010. GCN5-dependent acetylation of HIV-
1125 1 integrase enhances viral integration. *Retrovirology* **7**:18.
- 1126 80. **Topper M, Luo Y, Zhadina M, Mohammed K, Smith L, Muesing MA.** 2007.
1127 Posttranslational acetylation of the human immunodeficiency virus type 1 integrase carboxyl-
1128 terminal domain is dispensable for viral replication. *J Virol* **81**:3012–3017.
- 1129 81. **Zufferey R, Nagy D, Mandel RJ, Naldini L, Trono D.** 1997. Multiply attenuated lentiviral
1130 vector achieves efficient gene delivery in vivo. *Nat Biotechnol* **15**:871–875.
- 1131 82. **Rossillo P, Winter F, Simon-Loriere E, Gallois-Montbrun S, Negroni M.** 2012.
1132 Retroevolution: HIV-driven evolution of cellular genes and improvement of anticancer drug
1133 activation. *PLoS Genet* **8**:e1002904.
- 1134 83. **Adachi A, Gendelman HE, Koenig S, Folks T, Willey R, Rabson A, Martin MA.** 1986.
1135 Production of acquired immunodeficiency syndrome-associated retrovirus in human and
1136 nonhuman cells transfected with an infectious molecular clone. *J Virol* **59**:284–291.
- 1137 84. **Gasser R, Hamoudi M, Pellicciotta M, Zhou Z, Visdeloup C, Colin P, Braibant M, Lagane**
1138 **B, Negroni M.** 2016. Buffering deleterious polymorphisms in highly constrained parts of HIV-
1139 1 envelope by flexible regions. *Retrovirology* **13**:50.
- 1140 85. **Charneau P, Mirambeau G, Roux P, Paulous S, Buc H, Clavel F.** 1994. HIV-1 reverse
1141 transcription. A termination step at the center of the genome. *J Mol Biol* **241**:651–662.
- 1142 86. **Rosen CA, Sodroski JG, Campbell K, Haseltine WA.** 1986. Construction of recombinant
1143 murine retroviruses that express the human T-cell leukemia virus type II and human T-cell
1144 lymphotropic virus type III trans activator genes. *J Virol* **57**:379–384.
- 1145 87. **Foley Ge, Lazarus H, Farber S, Uzman Bg, Boone Ba, Mccarthy Re.** 1965. Continuous
1146 culture of human lymphoblasts from peripheral blood of a child with acute leukemia. *Cancer*
1147 **18**:522–529.
- 1148 88. **Nara PL, Fischinger PJ.** 1988. Quantitative infectivity assay for HIV-1 and-2. *Nature*
1149 **332**:469–470.
- 1150 89. **Nara PL, Hatch WC, Dunlop NM, Robey WG, Arthur LO, Gonda MA, Fischinger PJ.**
1151 1987. Simple, rapid, quantitative, syncytium-forming microassay for the detection of human
1152 immunodeficiency virus neutralizing antibody. *AIDS Res Hum Retroviruses* **3**:283–302.
- 1153 90. **Naldini L, Blömer U, Gallay P, Ory D, Mulligan R, Gage FH, Verma IM, Trono D.** 1996. In
1154 vivo gene delivery and stable transduction of nondividing cells by a lentiviral vector. *Science*
1155 **272**:263–267.

- 1156 91. **Vozzolo L, Loh B, Gane PJ, Tribak M, Zhou L, Anderson I, Nyakatura E, Jenner RG,**
1157 **Selwood D, Fassati A.** 2010. Gyrase B inhibitor impairs HIV-1 replication by targeting Hsp90
1158 and the capsid protein. *J Biol Chem* **285**:39314–39328.
- 1159 92. **De Iaco A, Santoni F, Vannier A, Guipponi M, Antonarakis S, Luban J.** 2013. TNPO3
1160 protects HIV-1 replication from CPSF6-mediated capsid stabilization in the host cell
1161 cytoplasm. *Retrovirology* **10**:20.
- 1162 93. **Sarzotti-Kelsoe M, Bailer RT, Turk E, Lin C-L, Bilka M, Greene KM, Gao H, Todd CA,**
1163 **Ozaki DA, Seaman MS, Mascola JR, Montefiori DC.** 2014. Optimization and validation of
1164 the TZM-bl assay for standardized assessments of neutralizing antibodies against HIV-1. *J*
1165 *Immunol Methods* **409**:131–146.
- 1166 94. **Lederle A, Su B, Holl V, Penichon J, Schmidt S, Decoville T, Laumond G, Moog C.**
1167 2014. Neutralizing antibodies inhibit HIV-1 infection of plasmacytoid dendritic cells by an
1168 FcγRIIIa independent mechanism and do not diminish cytokines production. *Sci Rep* **4**:5845.
- 1169 95. **Kabsch W.** 2010. Integration, scaling, space-group assignment and post-refinement. *Acta*
1170 *Crystallogr D Biol Crystallogr* **66**:133–144.
- 1171 96. **Kabsch W.** 2010. XDS. *Acta Crystallogr D Biol Crystallogr* **66**:125–132.
- 1172 97. **McCoy AJ, Grosse-Kunstleve RW, Adams PD, Winn MD, Storoni LC, Read RJ.** 2007.
1173 Phaser crystallographic software. *J Appl Crystallogr* **40**:658–674.
- 1174 98. **Adams PD, Afonine PV, Bunkóczi G, Chen VB, Davis IW, Echols N, Headd JJ, Hung L-**
1175 **W, Kapral GJ, Grosse-Kunstleve RW, McCoy AJ, Moriarty NW, Oeffner R, Read RJ,**
1176 **Richardson DC, Richardson JS, Terwilliger TC, Zwart PH.** 2010. PHENIX: a
1177 comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr*
1178 *D Biol Crystallogr* **66**:213–221.
- 1179 99. **Eijkelenboom AP, Sprangers R, Hård K, Puras Lutzke RA, Plasterk RH, Boelens R,**
1180 **Kaptein R.** 1999. Refined solution structure of the C-terminal DNA-binding domain of human
1181 immunovirus-1 integrase. *Proteins* **36**:556–564.
- 1182 100. **Terwilliger TC, Grosse-Kunstleve RW, Afonine PV, Moriarty NW, Adams PD, Read RJ,**
1183 **Zwart PH, Hung L-W.** 2008. Iterative-build OMIT maps: map improvement by iterative model
1184 building and refinement without model bias. *Acta Crystallogr D Biol Crystallogr* **64**:515–524.
- 1185 101. **Terwilliger TC, Grosse-Kunstleve RW, Afonine PV, Moriarty NW, Zwart PH, Hung L-W,**
1186 **Read RJ, Adams PD.** 2008. Iterative model building, structure refinement and density
1187 modification with the PHENIX AutoBuild wizard. *Acta Crystallogr D Biol Crystallogr* **64**:61–69.
- 1188 102. **Afonine PV, Grosse-Kunstleve RW, Echols N, Headd JJ, Moriarty NW, Mustyakimov M,**
1189 **Terwilliger TC, Urzhumtsev A, Zwart PH, Adams PD.** 2012. Towards automated
1190 crystallographic structure refinement with phenix.refine. *Acta Crystallogr D Biol Crystallogr*
1191 **68**:352–367.
- 1192 103. **Emsley P, Lohkamp B, Scott WG, Cowtan K.** 2010. Features and development of Coot.
1193 *Acta Crystallogr D Biol Crystallogr* **66**:486–501.
- 1194 104. **Pei J, Kim B-H, Grishin NV.** 2008. PROMALS3D: a tool for multiple protein sequence and
1195 structure alignments. *Nucleic Acids Res* **36**:2295–2300.
- 1196 105. **Krissinel E, Henrick K.** 2004. Secondary-structure matching (SSM), a new tool for fast
1197 protein structure alignment in three dimensions. *Acta Crystallogr D Biol Crystallogr* **60**:2256–
1198 2268.

- 1199 106. **Hough MA, Wilson KS.** 2018. From crystal to structure with CCP4. *Acta Crystallogr D Struct*
1200 *Biol* **74**:67–67.
- 1201 107. **Robert X, Gouet P.** 2014. Deciphering key features in protein structures with the new
1202 ENDscript server. *Nucleic Acids Res* **42**:W320–4.
- 1203 108. **Sarkar S, Witham S, Zhang J, Zhenirovskyy M, Rocchia W, Alexov E.** 2013. DelPhi Web
1204 Server: A comprehensive online suite for electrostatic calculations of biological
1205 macromolecules and their complexes. *Commun Comput Phys* **13**:269–284.
- 1206 109. **Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE.**
1207 2004. UCSF Chimera: a visualization system for exploratory research and analysis. *J*
1208 *Comput Chem* **25**:1605–1612.

Table 1. Efficiency of nuclear import and 3' processing ^a.

	wt IN A (NKNK)	IN A D116A	NQNQ	NQNK	NKNQ
Observed levels of integration 1 (relative to wt IN A) <i>Values from Figure 5</i>	1	0.00 ± 0.00	0.00 ± 0.00	0.03 ± 0.02	0.24 ± 0.09
Theoretical levels of 2LTRc 2 (relative to IN D116A) <i>see Materials and Methods</i>	0.20	1	1.00 ± 0.00	0.98 ± 0.02	0.81 ± 0.07
Observed levels of 2LTRc 3 (relative to IN D116A) <i>Values from Figure 7A</i>	0.20	1	0.33 ± 0.08	0.30 ± 0.09	0.28 ± 0.11
Efficiency of nuclear import 4 (relative to IN D116A) <i>Ratio values line 3 / values line 2</i>	1.00	1	0.33 ± 0.08	0.31 ± 0.09	0.35 ± 0.14
Ratio of PJ/2LTRc 5 (relative to IN D116A) <i>Values from figure 7B</i>	0.54 ± 0.02	1	0.93 ± 0.18	0.87 ± 0.19	0.66 ± 0.17
Decrease of PJ/2LTRc 6 (relative to IN D116A) <i>= 1-values in line 5</i>	0.46 ± 0.10	0	0.07 ± 0.01	0.13 ± 0.03	0.34 ± 0.09
Efficiency of 3' processing 7 (relative to wt IN A) <i>= values in line 6 / 0.46</i>	1.00 ± 0.22	0	0.15 ± 0.03	0.28 ± 0.06	0.74 ± 0.19
Expected levels of integration 8 (relative to wt IN A) <i>Product of values in lines 4 and 7</i>	1.00 ± 0.22	0	0.05 ± 0.01	0.09 ± 0.03	0.26 ± 0.11

^a SD values (calculated as described in Materials and Methods) are given in grey.

Table 2: Primers and probes used in the qPCR assay

duplex	target	primers/probes	sequence (5'-3')	fluorophore
quantification	U5Psi	U5Psi-forward	GTGACTCTGGTAACTAGAGA	-
	U5Psi	U5Psi-probe	CGCTTTCAAGTCCCTGTTCCGGG	FAM
	U5Psi	U5Psi-reverse	GAGAGCTCCTCTCCTTTC	-
	RU5	RU5-forward	CAGATCTGAGCCTGGGAG	-
	RU5	RU5-probe	AAGCAGTGGGTTCCCTAGTTAGCC	HEX
	RU5	RU5-reverse	GGCACACACTACTTGAAGC	-
normalisation	ACTB	IDT pre-designed assay, Hs.PT.56a.40703009.g/exon 6		HEX

Table 3: Primers used for the Alu qPCR assay

stage	target	primers/probes	sequence (5'-3')	fluorophore
1 st PCR	ALU-LTR	Alu forward	TGCTGGGATTACAGGCGTGAG	-
	ALU-LTR	Psi reverse	GTCCTCTGGTTCCCTTC	-
2 nd qPCR	RU5	RU5-forward	CAGATCTGAGCCTGGGAG	-
	RU5	RU5-probe	AAGCAGTGGGTTCCCTAGTTAGCC	HEX
	RU5	RU5-reverse	GGCACACACTACTTGAAGC	-

Table 4: Primers and probes used for the detection of total 2LTRc and 2LTRc with perfect palindromic junction (PJ).

target	primers/probes	sequence (5'-3')	fluorophore
2LTRc	2LTR forward	CCCTTTTAGTCAGTGTGGAA	-
2LTRc	2LTR probe	TTCACTCCCAACGAAGACAAGATATCCTT	FAM
2LTRc	2LTR reverse	GTAGCCTTGTGTGGTAGA	-
PJ	2LTR PJ forward	TGTGGAAAAATCTCTAGCAGTAC	-
PJ	2LTR probe	TTCACTCCCAACGAAGACAAGATATCCTT	FAM
PJ	2LTR reverse	GTAGCCTTGTGTGGTAGA	-

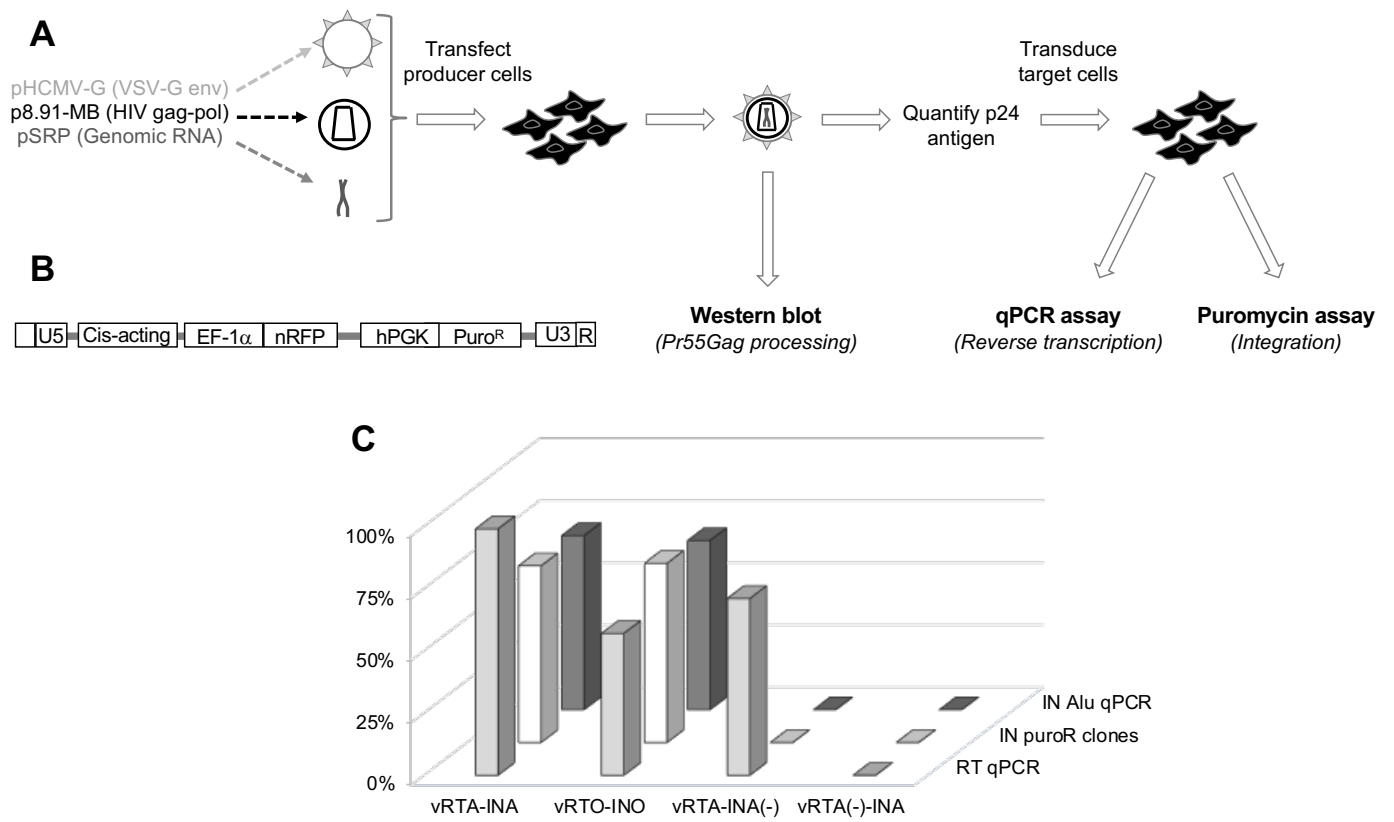


Figure 1

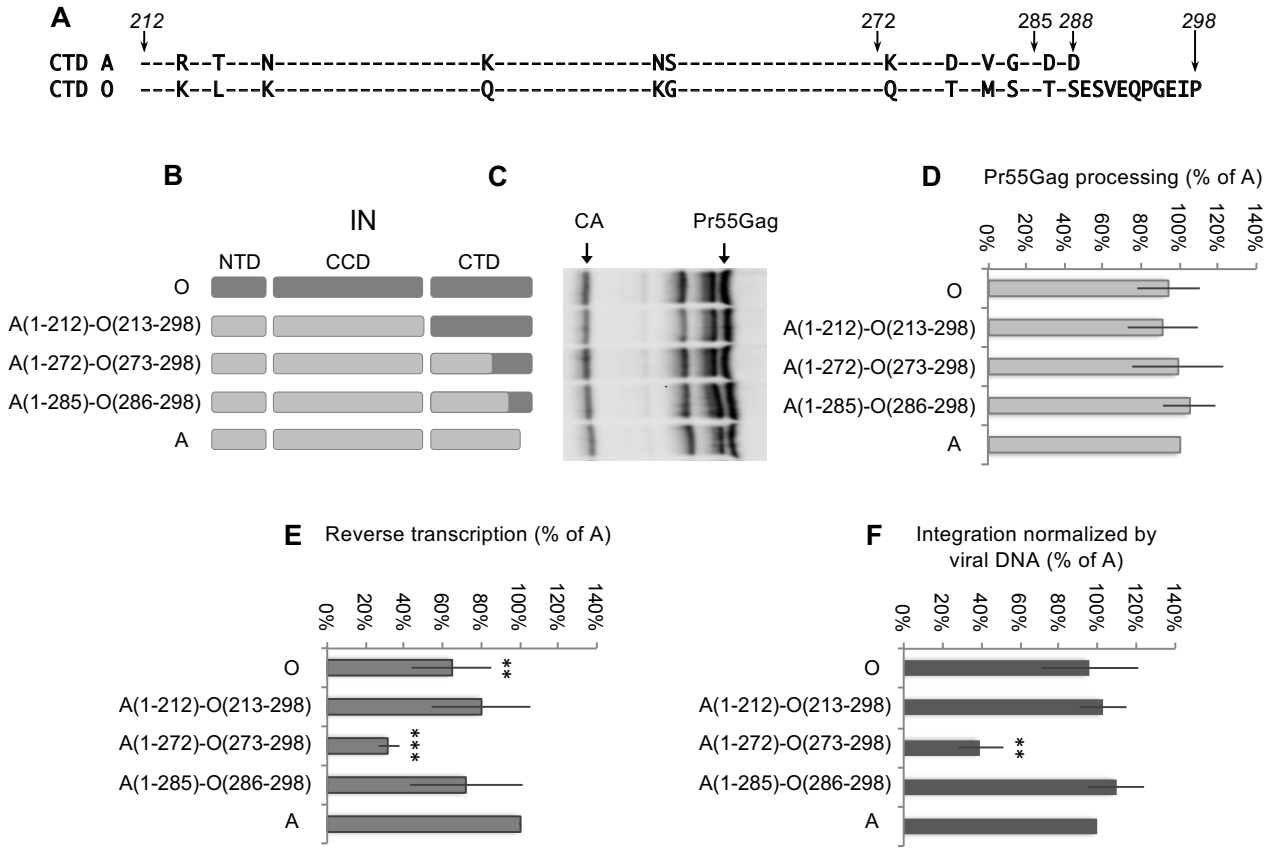


Figure 2

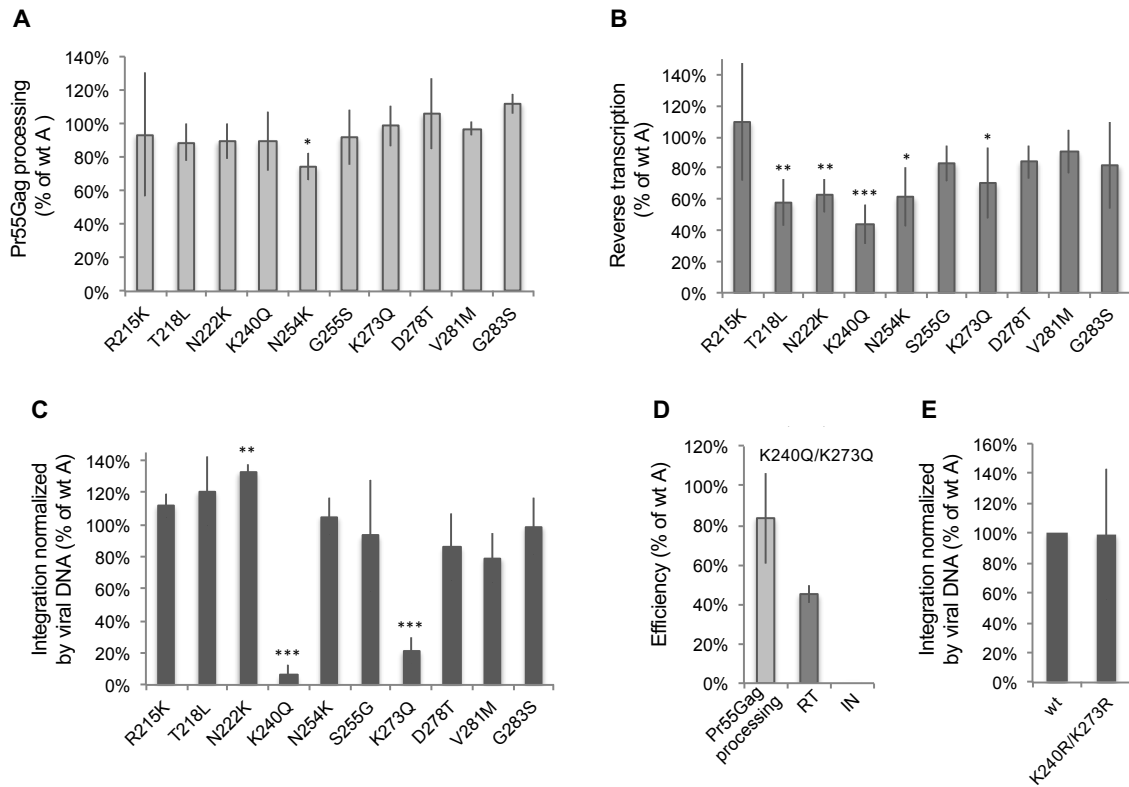


Figure 3

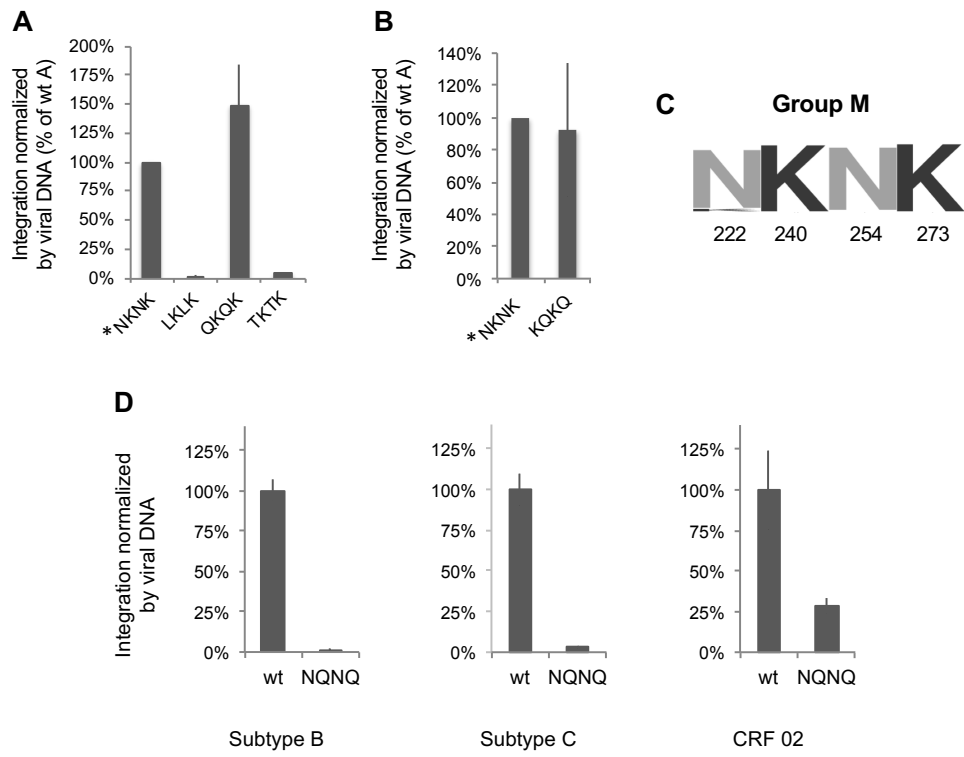


Figure 4

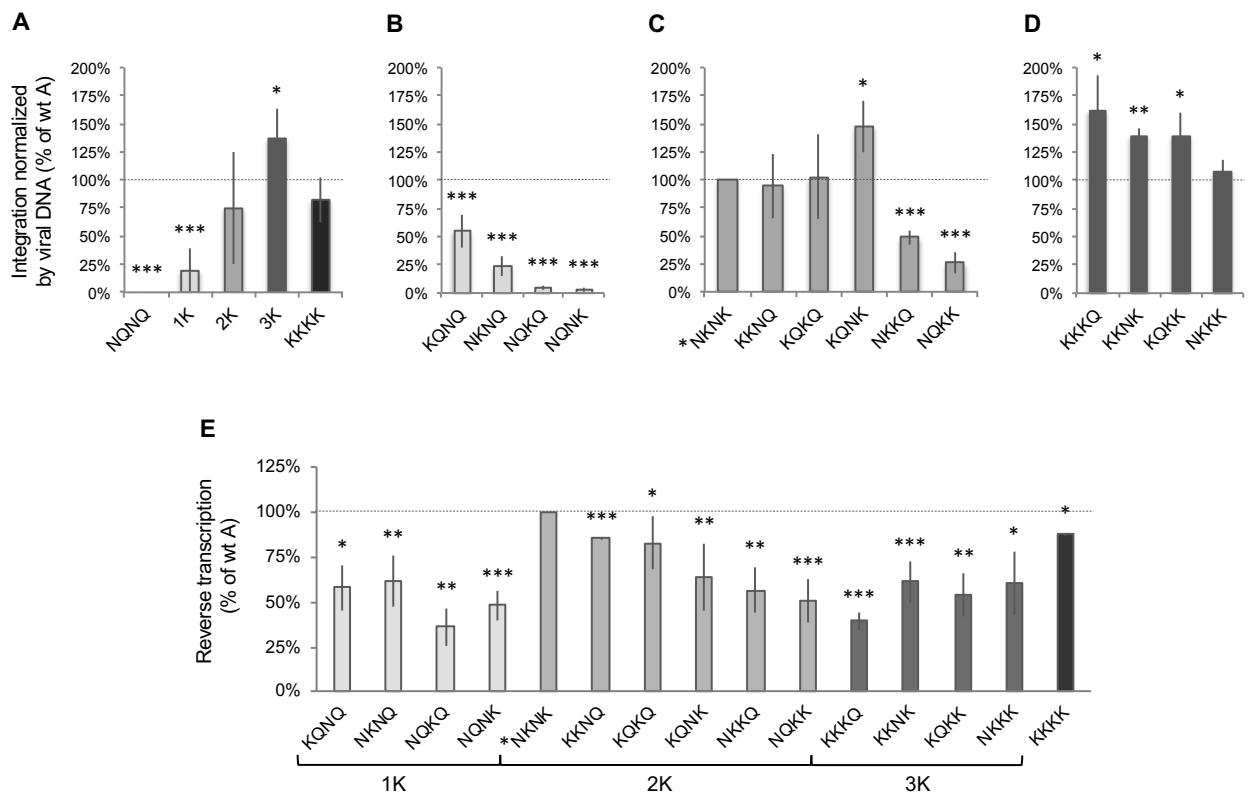


Figure 5

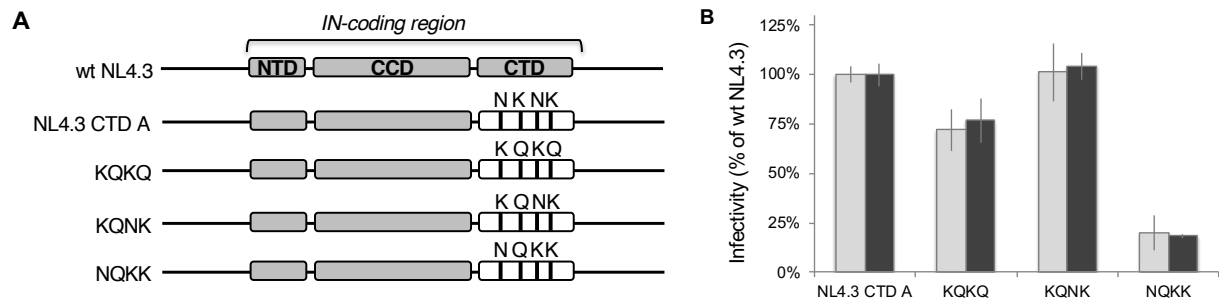
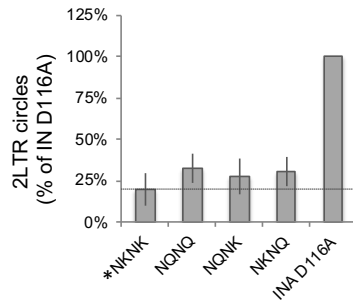


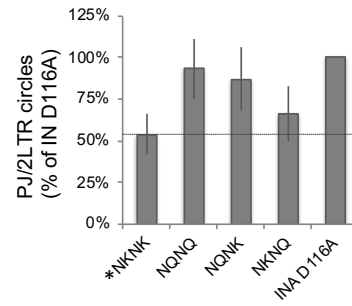
Figure 6

A

≠ to wt INA	na	*	ns	*	***
≠ to INA D116A	***	***	***	***	na
n	7	6	5	7	7

**B**

≠ to wt INA	na	***	***	ns	***
≠ to INA D116A	***	ns	ns	**	na
n	7	6	5	7	7

**Figure 7**

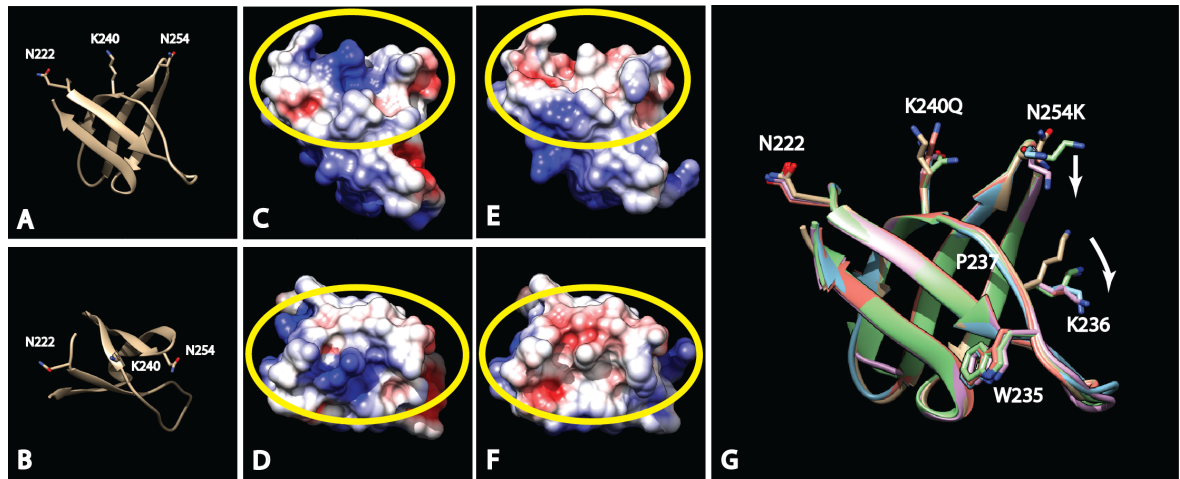


Figure 8

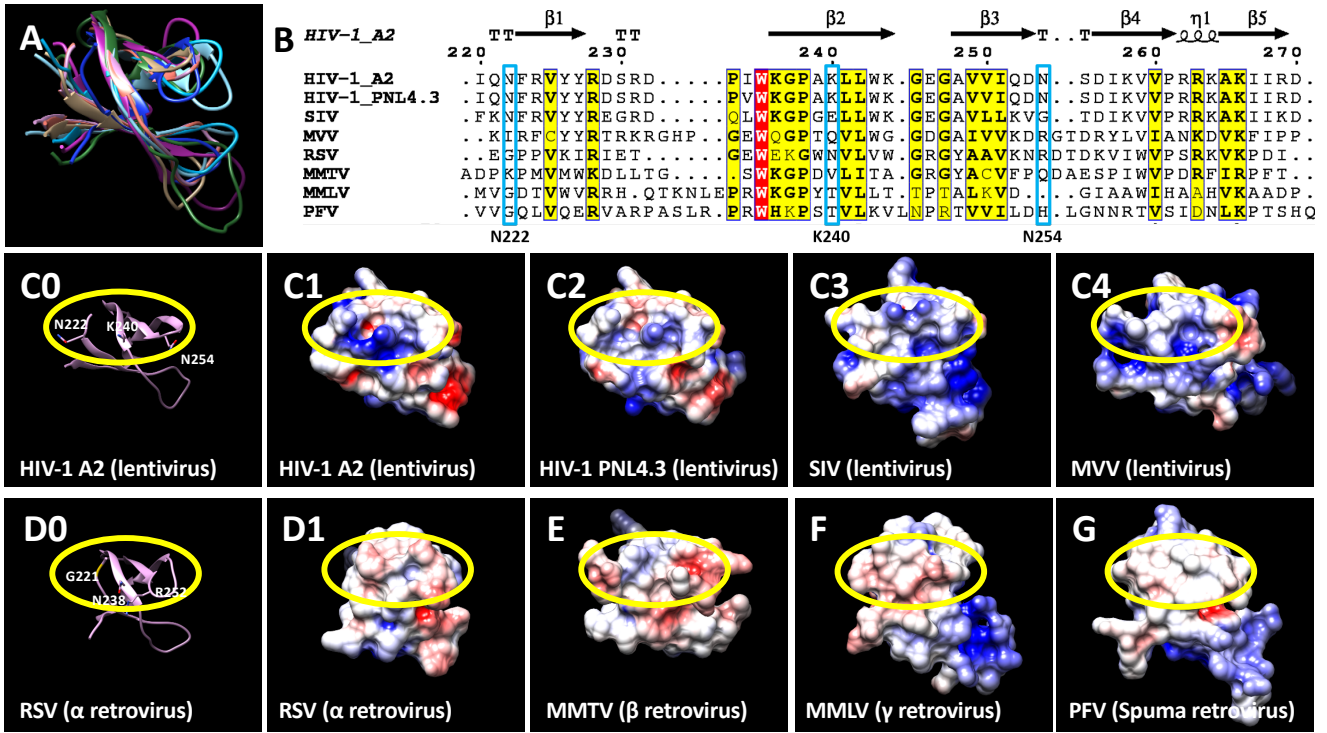


Figure 9

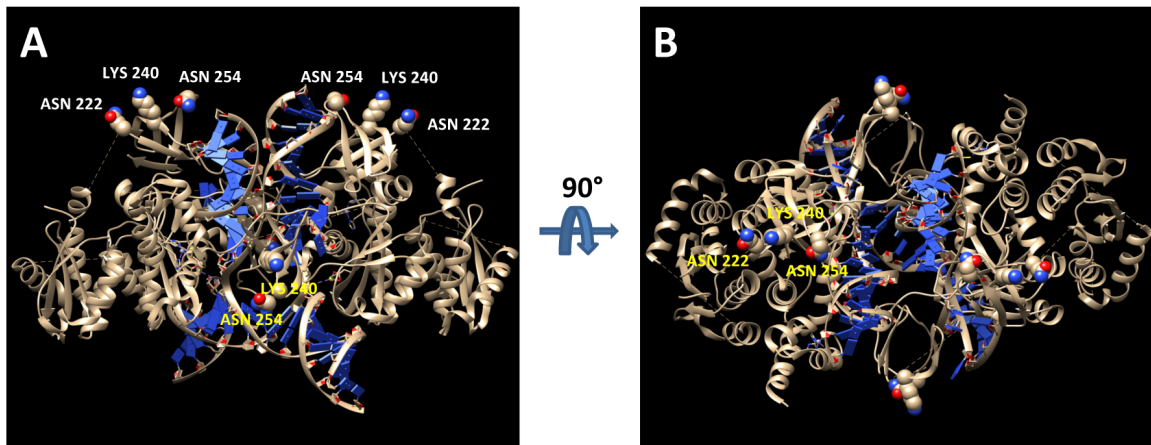


Figure 10